

APPLIED ANALYTIC MODELING

BA 706 FALL 2022

GROUP PROJECT

A FULL DATA MINING AND PREDICTIVE MODELING PROJECT AND REPORT ON

'EMPLOYEE TURNOVER DATASET'

(Gotten from kaggle.com)

IBUKUNOLUWA OLUKOKO-301259629

INTRODUCTION

Employee attrition happens when employees resign from a company at a rate higher than the company employs new hires. This can be a concern for a company as the cost implications from brain drain and reputational risk can become a challenge if this phenomenon is not checked.

We are assuming the problem statement is the high turnover of employees in different industries and our mission is to discover the reason for this high turnover rate and give recommendations on how to reduce the risk of losing good hands in the industries by reducing employee turnover.

Our data set, Employee Turnover (Kaggle.com) contains 16 variables and over 1000 observations. We will be predicting the variables that impact an employee's decision to resign from their place of employment.

The data was collected using questionnaires completed by the staff who quit. The employee database was also used to collate historical and present generic information about all employees.

There are several factors that could affect an employee's decision to leave their place of employment. Below is the dictionary showing the description of the variables we will be drawing our conclusions from in this project:

NAME	MODEL ROLE	MEASUREMENT LEVEL	DESCRIPTION
STAG	INPUT	INTERVAL	Experience(time) 0.39 - 179
EVENT	TARGET	BINARY	Employee turnover (0-NO or 1YES)
GENDER	INPUT	NOMINAL	Employee's gender (Female/Male)
AGE	INPUT	INTERVAL	Employee's Age
INDUSTRY	INPUT	NOMINAL	Employee's Industry
PROFESSION	INPUT	NOMINAL	Employee's profession
TRAFFIC	REJECTED	NOMINAL	From what pipeline did the employee come to the company (empjs, advert, rabrecNERab, youjs,referral,recNERab)
COACH	REJECTED	NOMINAL	Presence of training on probation (YES or NO)
HEAD_GENDER	INPUT	NOMINAL	Supervisor's gender (Female/Male)
GREYWAGE	INPUT	NOMINAL	Salary (White/Grey)
WAY	INPUT	NOMINAL	Employee's way of transportation

EXTRAVERSION	INPUT	INTERVAL	Extraversion score (1- 10)
INDEPEND	INPUT	INTERVAL	Independence score (1-10)
SELFCONTROL	INPUT	INTERVAL	Self -control score (1-10)
ANXIETY	INPUT	INTERVAL	Anxiety score (1-10)
NOVATOR	INPUT	INTERVAL	Innovator score (1-10)

Description of the Traffic and Greywage variables

TRAFFIC	<i>This describes the means by which the employee learned of the vacancy in the company and how they applied. For example, direct recommendation by a friend who is an employee or through a job site to mention a few. As we have chosen to reject the variable, we will not dwell much on defining the terms.</i>
GREYWAGE	Grey-wage: in Russia or Ukraine means that the employer pays just a tiny bit amount of salary above the white-wage White wage: minimum wage

DATA SETUP/EXPLORATION

We first imported the **employee dataset** and then made '**event**' our target variable as it is a binary variable indicating whether an employee resigned from the company.

We rejected '**traffic**' and '**coach**' because by narrowing down the scope of your analysis to factors directly tied to the workplace and job-related aspects, we can generate more actionable insights. Understanding why employees leave or stay within an organization is crucial for implementing effective retention strategies. By focusing on variables within the organization's control.

Results - Node: StatExplore Diagram: Employee

File Edit View Window

Interval Variables

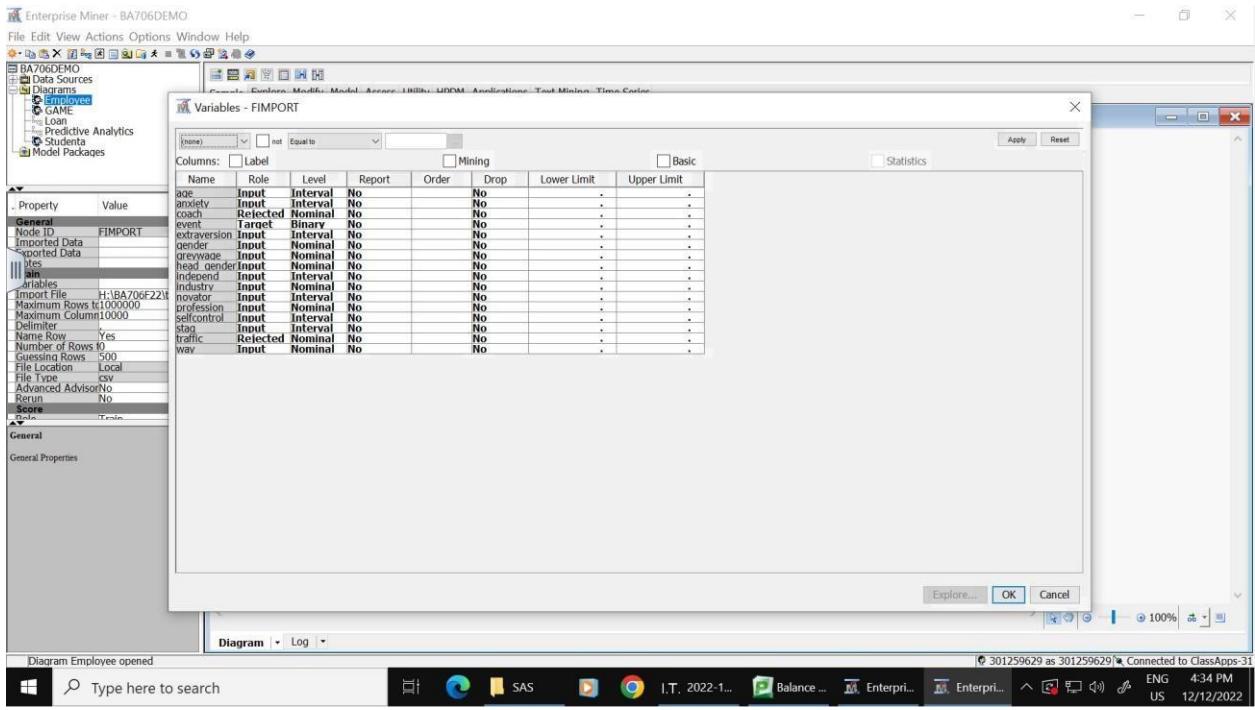
Data Role	Target	Target Level	Variable	Median	Missing	Non Missing	Minimum	Maximum	Mean	Standard Deviation	Skewness	Kurtosis	Role	Label	Scaled Mean Deviation	Maximum Deviation	Level Id
TRAIN	event	0	stad	30.225...	0	279	0.4928...	155.36	36.761	32.248...	1.8131	INPUT	stad	-0.0866	0.0584...	1	
TRAIN	event	1	stad	23.557...	0	285	0.3942...	160.06	34.442...	32.173...	1.8672...	INPUT	stad	-0.08641	0.0584...	2	
TRAIN	event	0	selfcon...	5.7	0	279	1	10.57878	1.9780...	0.0035...	-0.72717	INPUT	selfcon...	0.0302...	0.0295...	1	
TRAIN	event	1	selfcon...	5.7	0	285	1	10.54519	1.90378	0.1322...	-0.32953	INPUT	selfcon...	-0.0295...	0.0295...	2	
TRAIN	event	0	independ...	5.5	0	279	1	8.92053	1.9553...	-0.1926...	-0.10593	INPUT	independ...	-0.0287...	0.0277...	1	
TRAIN	event	1	independ...	5.5	0	285	1	10.55378	1.6874...	0.0380...	-0.16036	INPUT	independ...	0.0227...	0.0227...	2	
TRAIN	event	0	anxiety	5.6	0	279	1.7	10.579319	1.7066	0.0877...	-0.38356	INPUT	anxiety	0.0168...	0.0164...	1	
TRAIN	event	1	anxiety	5.6	0	285	1.7	9.456035...	1.7591...	0.2170...	-0.71277	INPUT	anxiety	-0.01647	0.0164...	2	
TRAIN	event	0	novator	8	0	279	1	10.59414	1.8340...	-0.2488...	-0.57474	INPUT	novator	0.0153...	0.0153...	1	
TRAIN	event	1	novator	8	0	285	1	10.59414	1.8340...	-0.2488...	-0.57474	INPUT	novator	0.0153...	0.0153...	2	
TRAIN	event	0	extrave...	5.4	0	279	1	10.55720...	1.9076...	-0.0799...	-0.17537	INPUT	extrave...	-0.00474	0.00464...	1	
TRAIN	event	1	extrave...	5.4	0	285	1	10.56245...	1.8409...	0.0590...	-0.52159	INPUT	extrave...	0.00464	0.00464...	2	
TRAIN	event	0	ace	30	0	279	19	54.31230...	6.4133...	0.6127...	-0.00539	INPUT	ace	0.0030...	0.00301...	1	
TRAIN	event	1	ace	30	0	285	18	54.31041...	7.0992...	0.48201...	-0.41285	INPUT	ace	-0.00301	0.00301...	2	

In the initial phase of statistical exploration we undertook a thorough examination of the dataset which yielded remarkable results - no missing values were observed across any variables considered. This attribute bolsters reliability and robustness for subsequent analysis since all available data is comprehensive enough to derive meaningful insights.

Furthermore an initial assessment determined that there was only minimal skewness present in the distribution. Skewness indicates asymmetry in the distribution with positive skewness suggesting a longer tail on the right side. In this particular case the observed skewness hinted at a minor deviation from perfect symmetry.

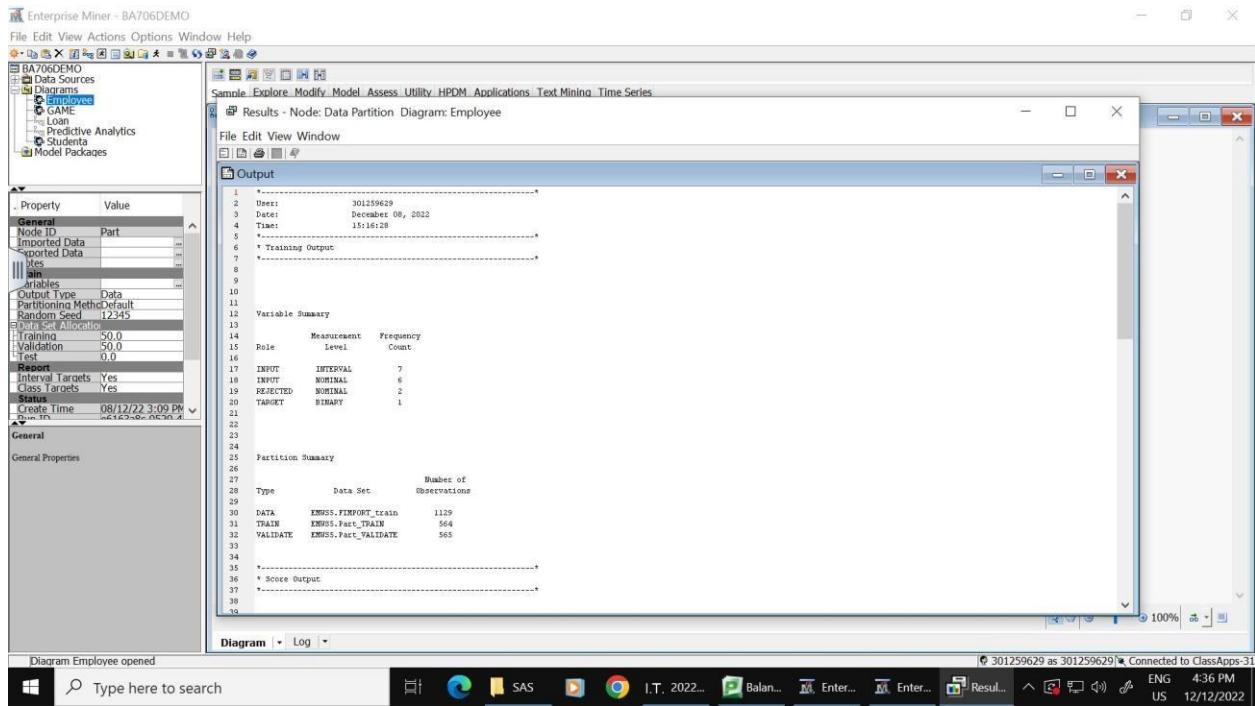
Although this skewness may have minimal impact on our current analysis. It is important to acknowledge its presence and make plans to address it during future stages of the project. Proper techniques for handling skewed data, such as logarithmic or power transformations, or utilizing specific statistical models catered to skewed distributions should be considered.

By recognizing and planning to address this slight skewness in subsequent stages of analysis we can ensure a more comprehensive and rigorous examination. Attending to these details will help guarantee accuracy in our results while fostering a deeper understanding of the factors contributing to employee attrition within our organization.



A data partition node was connected to the dataset with the below allocations:

Data Partition	Allocation	# of observations
Training	50%	564
Validation	50%	565
Test	0	0



DECISION TREES

The first set of models we used are decision trees. The trees are classified based on the method used , the number of branches, and the assessment measure. The following tree models were built:

1. Maximal Trees(Two and Three Branch)
2. Classification Trees(Two and Three Branch) 3. Probability Trees(Two and Three Branch)

TWO-BRANCH MAXIMAL TREE

The decision tree node was attached to the data partition node with a maximum of two branches indicated, ‘largest’ used as the method and ‘decision’ as the assessment measure. Below shows the outcome after this node was run with the specifications noted.

Enterprise Miner - BA706DEMO

File Edit View Actions Options Window Help

BA706DEMO

- Data Sources
- Diagrams
- Employee** (selected)
- GAME
- Loan
- Predictive Analytics
- Score Sheets
- Model Packages

Property Value

- Split Size: 1
- Search: No
- Use Dimensions: No
- Pruned: No
- MaxDepth: 5000
- MinSample: 20000
- Method: Largest
- Number of Leaves: 1
- Assessment Measure: Decision
- Assessment Fraction: 0.25
- Cross Validation: Perform Cross Validation: No
- Number of Folds: 10
- Number of Repeats: 1
- Seed: 12345
- Collaboration Based Imputation
- Number Single Var Imp: 0
- Bayesian Adjustment

General Properties

Diagram Log

Type here to search

SAS I.T. 2022-12-0... Balance f... Enterprise... Enterprise... Results...

ENG 2:43 PM 12/14/2022

Score Rankings Overlay: event

Cumulative Lift

Leaf Statistics

Treemap

Tree

Fit Statistics

Target	Target Label	Fit Statistics	Statistics	Train	Validation	Test
event	NOBS	Sum of Freq.	564	565		
event	MAX	Maximum A.	0.916697	0.916697		
event	SSE	Sum of Squ.	243.6498	276.9288		
event	ASE	Average Sq.	0.216002	0.24507		

Output

1 User: 301259629
 2 Date: December 14, 2022
 3 Time: 11:40:59
 4 *
 5 *
 6 * Training Output
 < *>

Results - Node: Max Decision Tree. Diagram: Employee

File Edit View Window

Tree

Node Id: 1 Statistic: Train Validation O: 45.47% 45.38% I: 54.52% 54.62% Count: 564 SED5

Node Id: 2 Statistic: Train Validation O: 50.40% 56.01% I: 41.60% 43.98% Count: 551 241

Node Id: 3 Statistic: Train Validation O: 34.74% 38.29% I: 65.25% 60.71% Count: 213 224

Node Id: 4 Statistic: Train Validation O: 62.64% 60.00% I: 37.36% 40.00% Count: 296 280

Node Id: 5 Statistic: Train Validation O: 34.55% 37.70% I: 65.45% 62.30% Count: 15 61

Node Id: 6 Statistic: Train Validation O: 54.84% 53.33% I: 45.16% 46.67% Count: 31 31

Node Id: 7 Statistic: Train Validation O: 31.32% 37.11% I: 68.68% 62.88% Count: 182 194

Node Id: 8 Statistic: Train Validation O: 23.08% 28.30% I: 76.14% 71.70% Count: 42 47

Node Id: 9 Statistic: Train Validation O: 69.75% 38.71% I: 30.25% 64.29% Count: 13 14

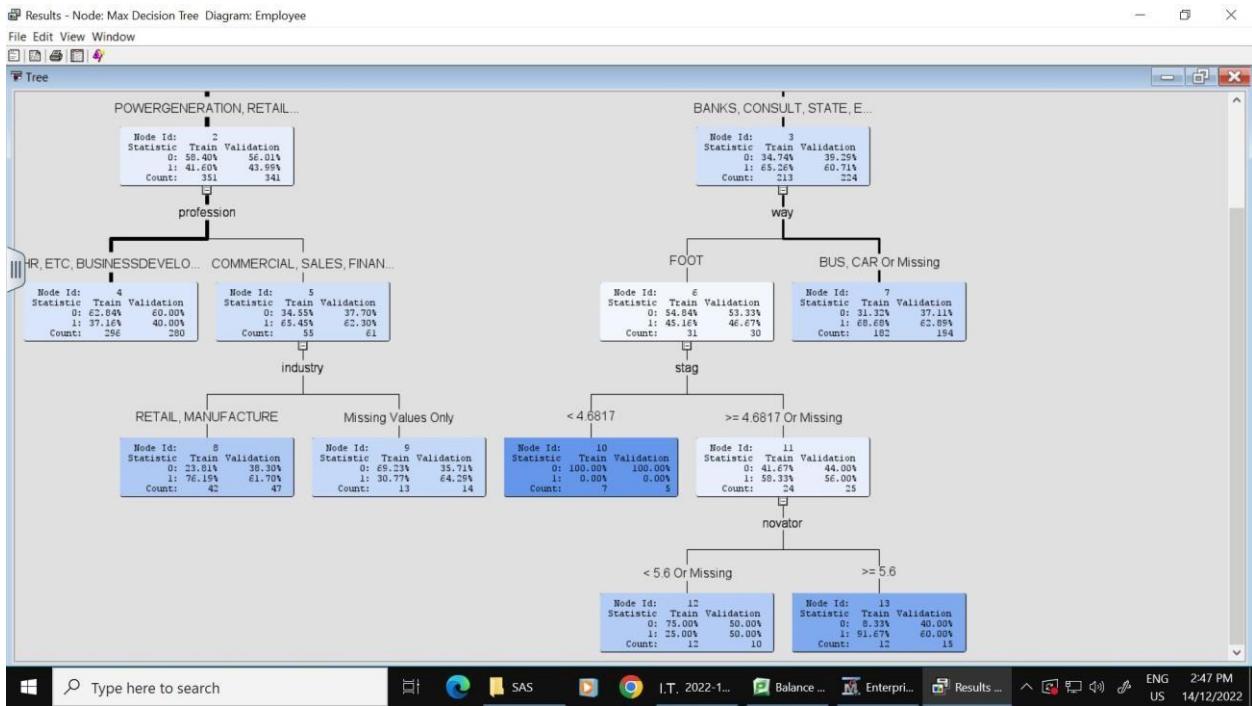
Node Id: 10 Statistic: Train Validation O: 100.00% 100.00% I: 0.00% 0.00% Count: 7 5

Node Id: 11 Statistic: Train Validation O: 44.00% 44.00% I: 56.33% 56.00% Count: 24 25

Type here to search

SAS I.T. 2022-1... Balance f... Enterprise... Enterprise... Results...

ENG 2:47 PM 12/14/2022



Results - Node: Max Decision Tree Diagram: Employee

File Edit View Window

Fit Statistics

Target	Target Label	Fit Statistics	Statistics Label	Train	Validation	Test
event		NOBS	Sum of Frequencies	564	565	
event		MISC	Misclassification Rate	0.328014	0.39292	
event		MAX	Maximum Absolute Error	0.916667	0.916667	
event		SSE	Sum of Squared Errors	243.6498	270.6698	
event		ASE	Average Squared Error	0.216002	0.24507	
event		RASE	Root Average Squared Error	0.46476	0.495045	
event		DIV	Divisor for ASE	1128	1130	
event		DFT	Total Degrees of Freedom	564		

```

Results - Node: Max Decision Tree Diagram: Employee
File Edit View Window
Output
57 -----
58
59
60
61 Variable Importance
62
63
64           Number of
65           Splitting      Ratio of
66 Variable   Label     Rules       Importance    Validation
67 Name        Importance      Importance    to Training
68 industry      2          1.0000      1.0000      1.0000
69 profession   1          0.6261      0.8761      1.3993
70 novatot     1          0.5306      0.0000      0.0000
71 stag         1          0.4413      0.6450      1.4616
72 way         1          0.3934      0.3951      1.0045
73
74
75
76 Tree Leaf Report
77
78           Training
79           Percent    Validation
80 Node      Observations   i   Observations   Validation
81 Id        Depth
82 4          2          296      0.37      280      0.40
83 7          2          182      0.69      194      0.63
84 8          3          42       0.76      47       0.62
85 9          3          13       0.31      14       0.64
86 12         4          12       0.25      10       0.50
87 13         4          12       0.92      15       0.60
88 10         3          7        0.00      5        0.00
89
90
91
92
93 Fit Statistics
94
95 Target-event Target Label=' '
96
97 Fit
98 Statistics Statistics Label      Train   Validation
99

```

The tree had seven leaves with five variable splits. The first split occurs at the Industry variable with the highest logworth when compared to the other variables selected in the model. The tree shows that 62.89% of employees in Banks, Consult and State who go to work in vehicles are likely to quit their jobs. 60% of employees with innovator score of greater than 5.6 who have worked greater than 4.7 months in the Bank, Consult, and State industries are likely to quit. 62% of Commercial, Sales, and Finance professionals in the power generation industry are likely to quit.

The Validation Average Square Error for the two-branch Maximal Tree is 0.24507.

THREE-BRANCH MAXIMAL TREE

The decision tree node was attached to the data partition node with a maximum of three branches indicated, ‘largest’ used as the method and ‘decision’ as the assessment measure. Below shows the outcome after this node was run with the specifications noted.

Enterprise Miner - BA706DEMO

File Edit View Actions Options Window Help

BA706DEMO

- Data Sources
- Diagrams
 - Employee**
 - GAME
 - Loan
 - Predictive Analytics
 - Students
- Model Packages

Results - Node: Max Decision Tree 3 maxB Diagram: Employee

Leaf Statistics

Node	Training Percent 1	Validation Percent 1
1	~0.55	~0.15
2	~0.45	~0.25
3	~0.35	~0.05
4	~0.50	~0.20
5	~0.25	~0.10
6	~0.15	~0.05
7	~0.10	~0.10
8	~0.05	~0.20
9	~0.05	~0.30

Score Rankings Overlay: event

Cumulative Lift

Fit Statistics

Target	Target Label	Fit Statistics	Statistics Label	Train	Validation	Test
event	NOBS	Sum of Fr...	564	565		
event	MISC	Misclassifi...	0.324468	0.40177		
event	MAX	Maximum ...	0.857143	0.857143		

Treemap

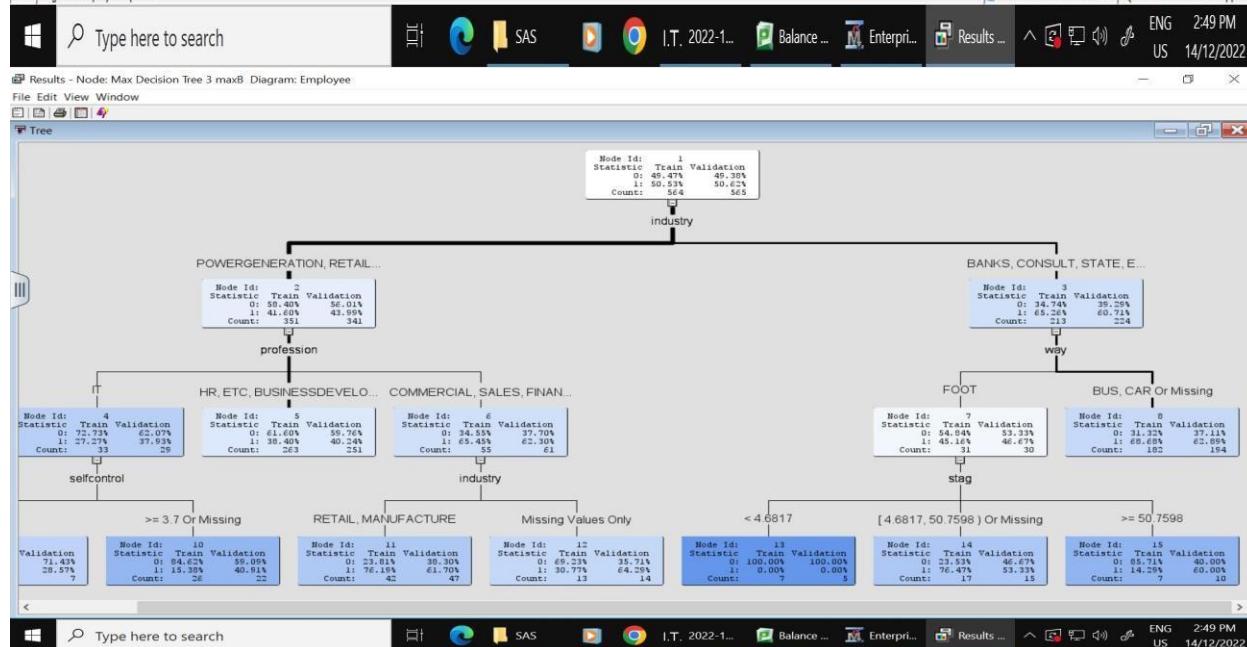
Output

```

1 -----
2 User:          301259629
3 Date:         December 14, 2022
4 Time:        11:52:31
5 -----

```

Diagram | Log | 301259629 as 301259629 Connected to ClassApps-25



Results - Node: Max Decision Tree 3 maxB Diagram: Employee

File Edit View Window

Fit Statistics

Target	Target Label	Fit Statistics	Statistics Label	Train	Validation	Test
event		N OBS	Sum of Frequencies	564	565	
event		MISC	Misclassification Rate	0.324468	0.40177	
event		MAX	Maximum Absolute Error	0.857143	0.857143	
event		SSE	Sum of Squared Errors	240.9574	283.4896	
event		ASE	Average Squared Error	0.213615	0.250876	
event		RASE	Root Average Squared Error	0.462185	0.500875	
event		DIV	Divisor for ASE	1128	1130	
event		DFT	Total Degrees of Freedom	564		

Results - Node: Max Decision Tree 3 maxB Diagram: Employee

File Edit View Window

Output

```

61 Variable Importance
62
63
64          Number of
65          Splitting
66 Variable      Rules   Importance       Validation
67 Name        Label    Importance       Importance
68 industry      2      1.0000      1.0000      1.0000
69 profession    1      0.6560      0.8436      1.2859
70 stag         1      0.6302      0.0000      0.0000
71 selfcontrol   1      0.4277      0.0000      0.0000
72 way          1      0.3934      0.3951      1.0045
73
74
75 Tree Leaf Report
76
77
78          Training
79          Percent
80 Node      Depth Observations   1   Observations   Validation
81 Id
82 5      2      263      0.38      251      0.40
83 8      2      182      0.69      194      0.63
84 11     3      42       0.76      47       0.62
85 10     3      26       0.15      22       0.41
86 14     3      17       0.76      15       0.53
87 9      3      13       0.31      14       0.64
88 9      3      7       0.71      7       0.29
89 13     3      7       0.00      5       0.00
90 15     3      7       0.14      10      0.60
91
92
93
94
95 Fit Statistics
96
97 Target=event Target Label=' '
98
99
100 Statistics Statistics Label      Train   Validation
101 _N OBS_ Sum of Frequencies      564.00      565.00
102 _MISC_ Misclassification Rate  0.32          0.40

```

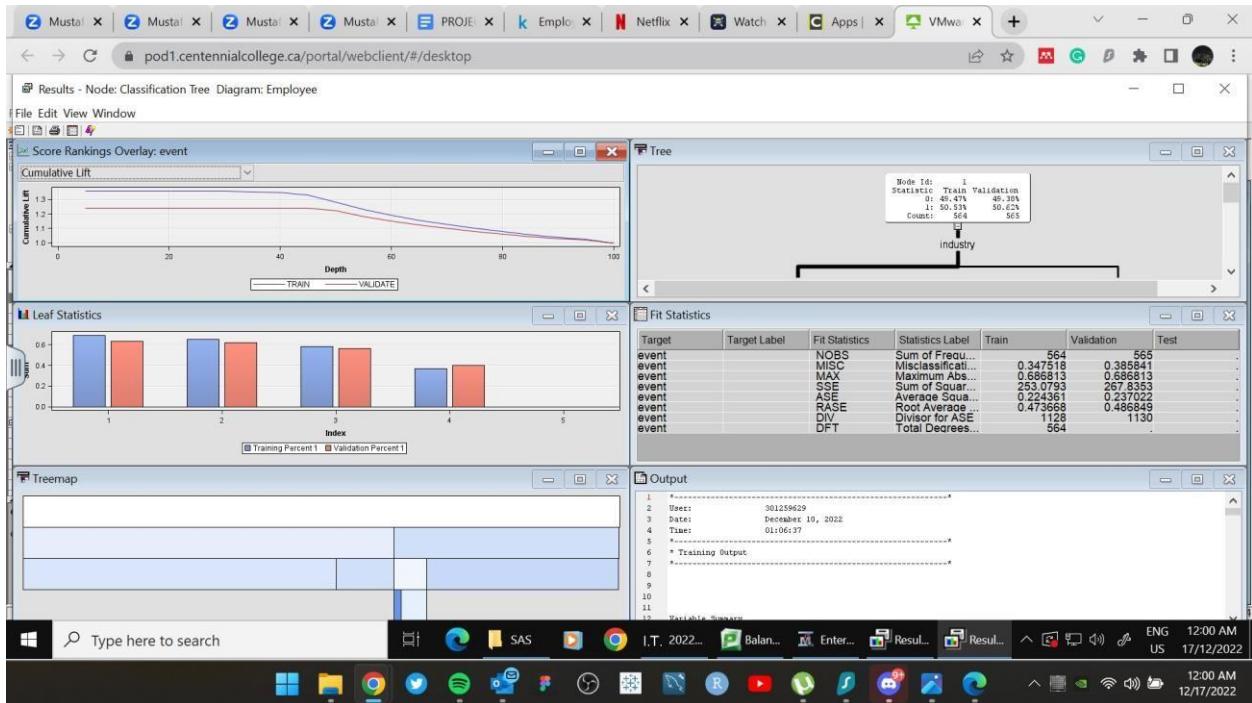
Type here to search

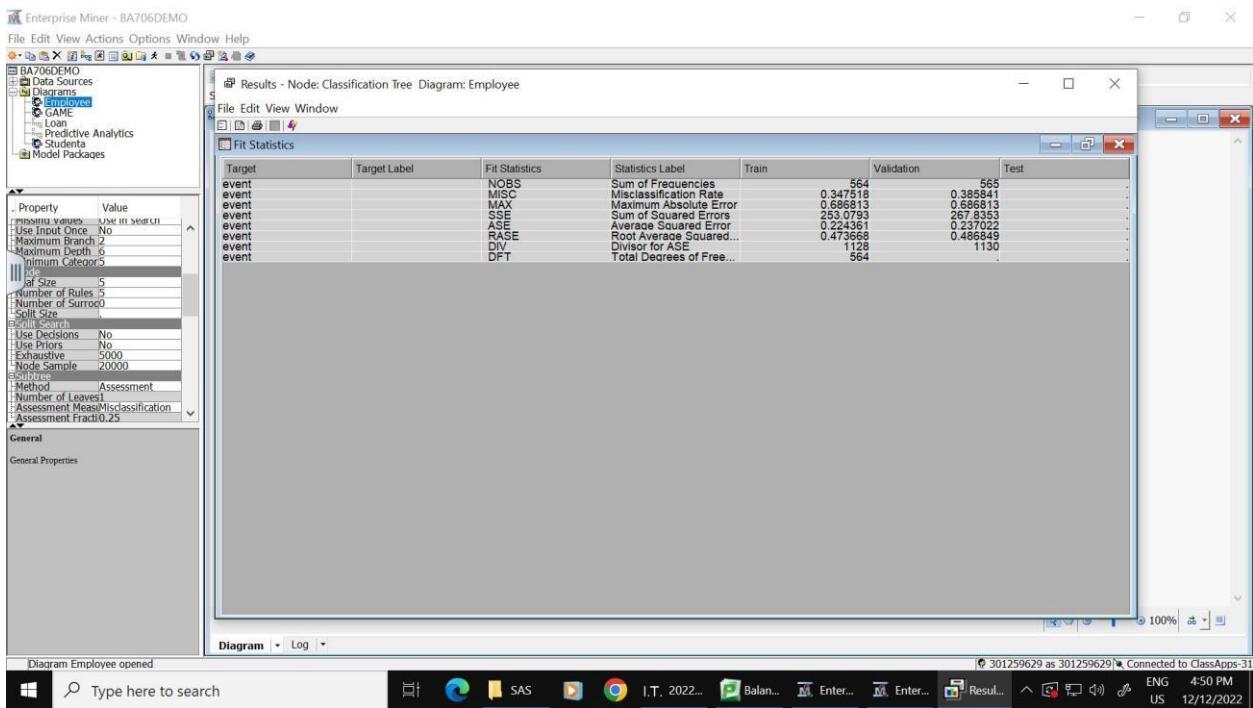
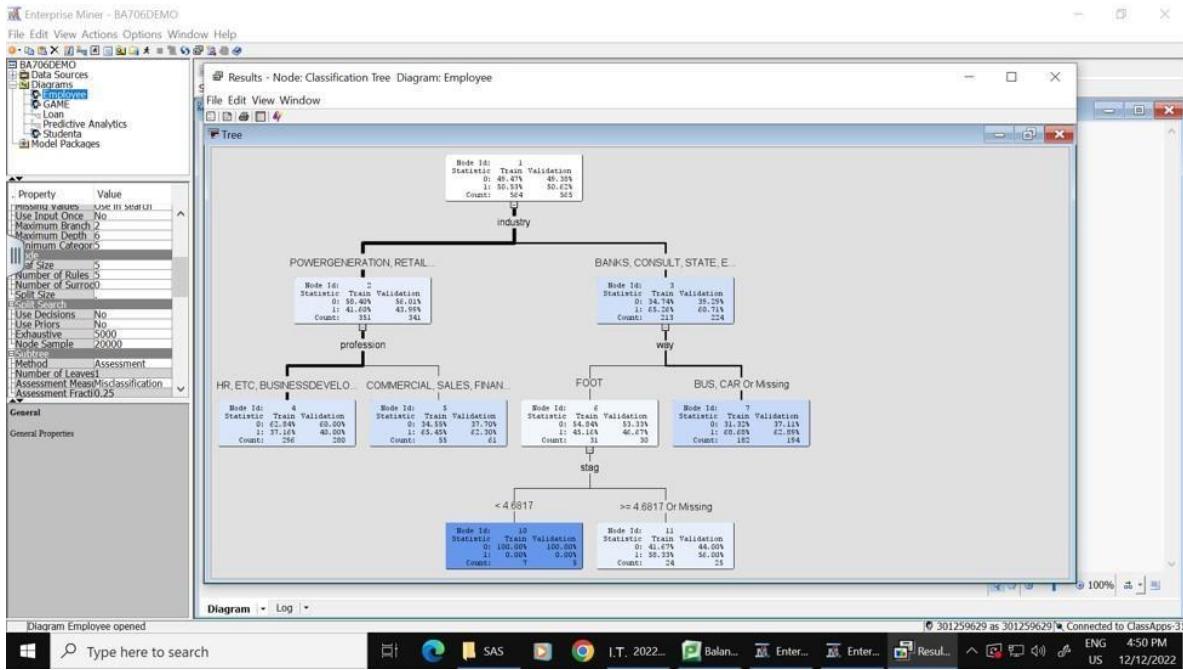
The tree had nine leaves with five variable splits. The first split occurs at the Industry variable with the highest logworth when compared to the other variables selected in the model. The tree shows that 62.89% of employees in Banks, Consult and State who go to work in vehicles are likely to quit their jobs. 60% of employees who have worked greater than 51 months in the Bank, Consult, and State industries and walk to the office are likely to quit. 62.3% of Commercial, Sales, and Finance professionals in the power generation industry are likely to quit.

The Validation Average Square Error for the three-branch Maximal Tree is 0.250876. We observed that the Two Branch Maximal Tree is a better model than the Three Branch Maximal tree, as this has a higher Validation Average Square Error.

TWO-BRANCH CLASSIFICATION TREE

The decision tree node was attached to the data partition node with a maximum of two branches indicated, ‘assessment’ was used for the method and ‘misclassification’ as the assessment measure. Below shows the outcome after this node was run with the specifications noted.





The screenshot shows the SAS Enterprise Miner interface with a classification tree node for 'Employee'. The tree has five leaves with four variable splits. The first split occurs at the Industry variable with the highest logworth when compared to the other variables selected in the model. The tree shows that 62.89% of employees in Banks, Consult and State industries who go to work in vehicles are likely to quit their jobs. 56% of employees who have worked greater than 4.7 months in the Bank, Consult, and State industries and walk to the office are likely to quit. 62.3% of Commercial, Sales, and Finance professionals in the power generation industry are also likely to quit.

```

42 PREDICTED_P_event1 Predicted: event=1
43 RESIDUAL_P_event1 Residual: event=1
44 PREDICTED_P_event0 Predicted: event=0
45 RESIDUAL_P_event0 Residual: event=0
46 FROM _P_event From: event
47 INTO _I_event Into: event
48
49
50 -----
51 * Score Output
52 -----
53
54 -----
55 * Report Output
56 -----
57
58
59
60
61 Variable Importance
62
63
64 Variable Number of Splitting Rules Validation Importance Ratio of Validation to Training Importance
65 Name label
66
67
68 industry 1 1.0000 1.0000 1.0000
69 profession 1 0.7073 0.6761 1.2388
70 state 1 0.4984 0.6450 1.2939
71 way 1 0.4443 0.3931 0.8892
72
73
74
75 Tree Leaf Report
76
77 Node Training Percent Validation Observations Validation Percent
78 Id Depth Observations 1
79
80
81 4 2 296 0.37 280 0.40
82 7 2 102 0.69 194 0.63
83 8 2 85 0.61 61 0.61

```

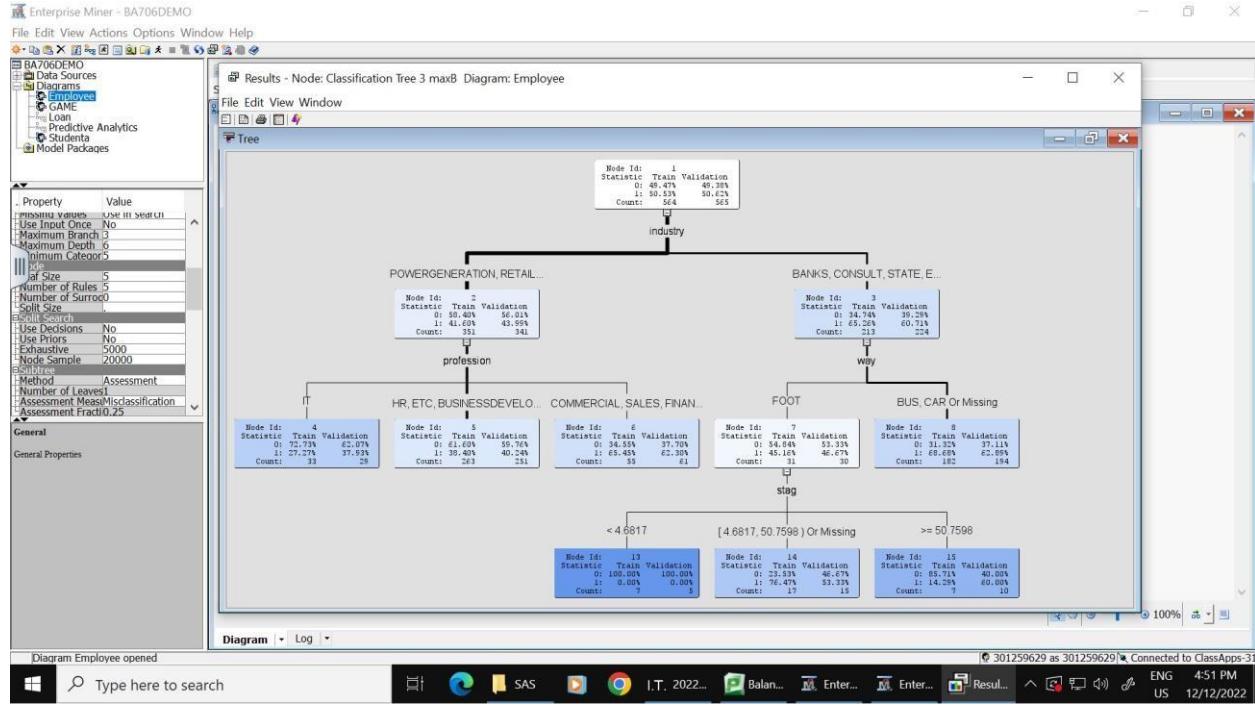
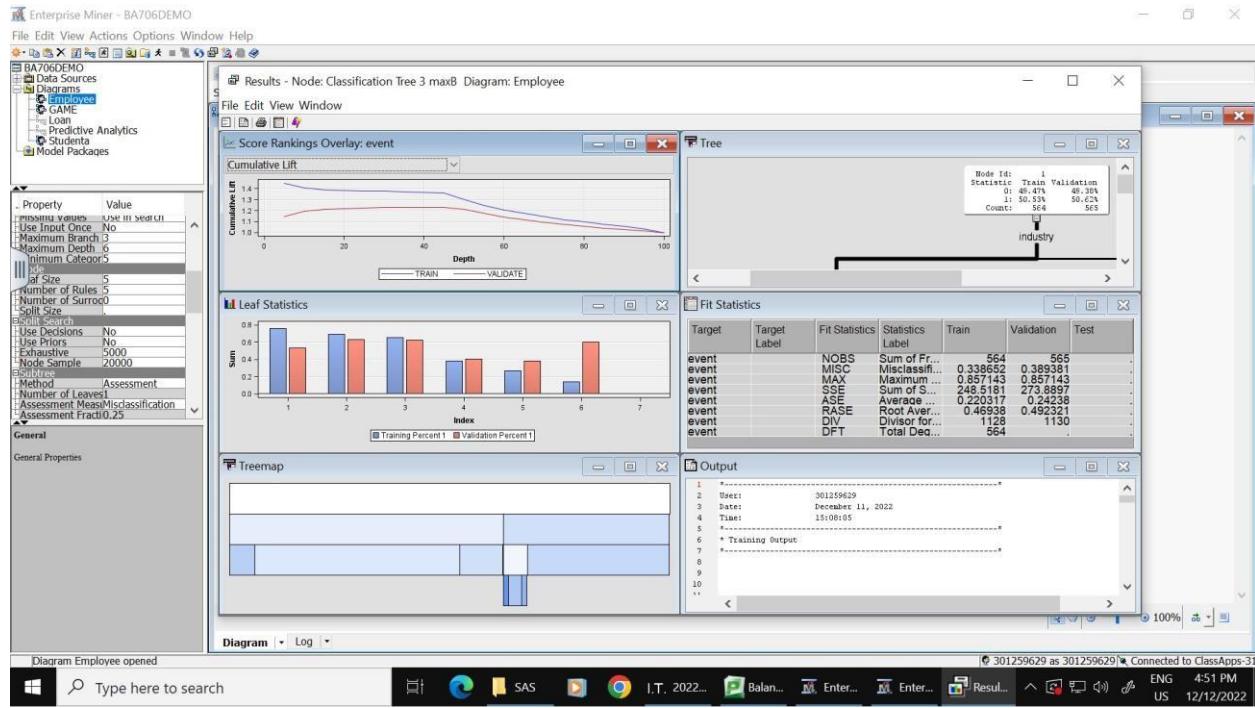
Diagram Employee opened

The tree has five leaves with four variable splits. This reduction is a result of pruning using the misclassification measure. The first split occurs at the Industry variable with the highest logworth when compared to the other variables selected in the model. The tree shows that 62.89% of employees in Banks, Consult and State industries who go to work in vehicles are likely to quit their jobs. 56% of employees who have worked greater than 4.7 months in the Bank, Consult, and State industries and walk to the office are likely to quit. 62.3% of Commercial, Sales, and Finance professionals in the power generation industry are also likely to quit.

The Validation Average Square Error for the two-branch classification Tree is 0.237022.

THREE-BRANCH CLASSIFICATION TREE

The decision tree node was attached to the data partition node with a maximum of three branches indicated, ‘assessment’ used as the method and ‘misclassification’ as the assessment measure. Below shows the outcome after this node was run with the specifications noted.



The screenshot shows the SAS Enterprise Miner interface with the following details:

- Project Tree:** BA706DEMO > Diagrams > Employee.
- Properties Panel:** Shows settings like 'Property' and 'Value' for various parameters such as 'Input Once' (No), 'Maximum Depth' (5), 'Number of Rules' (5), and 'Assessment' (Assessment).
- Output Window:** Displays the generated R code for the classification tree. Key parts include:
 - PREDICTED_P_event0 Predicted: event=0
 - RESIDUAL_P_event0 Residual: event=0
 - FROM_P_event From: event
 - INFO_I_event Info: event
 - Score Output
 - Report Output
- Variable Importance:** A table showing the importance of variables (industry, profession, sex, way) across training and validation datasets.
- Tree Leaf Report:** A table showing the distribution of observations across tree nodes (Node ID, Depth, Training Observations, Validation Observations, Percent).
- Status Bar:** Shows the diagram name (Employee), log status, and system information (Connected to ClassApps-31, 301259629 as 301259629, ENG, 4:52 PM, US, 12/12/2022).

The tree has seven leaves with four variable splits. The increased number of branches also increased the number of leaves when compared to the two branch model. The tree shows that 62.89% of employees in Banks, Consult and State who go to work in vehicles are likely to quit their jobs. 60% of employees who have worked greater than 51 months in the Bank, Consult, and State industries and walk to the office are likely to quit. 62.3% of Commercial, Sales, and Finance professionals in the power generation industry are likely to quit.

The Validation Average Square Error for the three-branch classification Tree is 0.24238. The Two Branch Classification Tree is a better model than the Three Branch Classification tree, as it has a lower Validation Average Square Error.

TWO BRANCH PROBABILITY TREE

The decision tree node was attached to the data partition node with a maximum branch of two, ‘assessment’ used as the method and ‘average squared error’ as the assessment measure. Below shows the outcome after this node was run with the specifications noted.

The screenshot displays the Enterprise Miner software interface with the following components:

- Left Sidebar:** Shows project navigation with nodes like BA706DEMO, Data Sources, Diagrams, and Predictive Analytics.
- Properties Panel:** Lists project properties such as Property, Value, Missing Values, Use in Search, Use Input Once, No, Maximum Depth, Minimum Category, Leaf Size, Number of Rules, Number of Surrogate Rules, Split Size, and Search Type.
- General Properties:** Section for general project settings.
- Diagram Employee opened:** Status bar at the bottom left.
- Top Menu:** File, Edit, View, Actions, Options, Window, Help.
- Score Rankings Overlay: event** (Cumulative Lift): A chart showing Cumulative Lift vs. Depth (TRAIN and VALIDATE) for an event.
- Tree:** A decision tree visualization with Node ID 1 (industry) and its statistics.
- Leaf Statistics:** Bar chart showing Sum vs. Index (Training Percent 1, Validation Percent 1).
- Fit Statistics:** Table showing Fit Statistics and Statistics for various events.
- Treemap:** A hierarchical treemap visualization.
- Output:** Text output window showing training and validation details.
- Bottom Navigation:** Diagram, Log, SAS, Google Chrome, I.T. 2022..., Balance Sheet, Enter..., Result..., and a search bar.

The screenshot shows the SAS Enterprise Miner interface with the following details:

- Title Bar:** Type here to search, SAS, I.T. 2022..., Balan..., Enter..., Enter..., Result..., ENG US 12/12/2022
- Left Panel (Project Tree):** BA706DEMO, Data Sources, Diagrams, Employee, GAME, Lasso, Predictive Analytics, Students, Model Packages.
- Left Panel (Properties):** Property Value, including: max/min values, Use in search, Use Input Once, No; Maximum Branch 2, Maximum Depth 5, minimum Categorical Size, minimum Categorical S, of Size 5, Number of Rules 5, Number of Surrogate 0, Split Size 1, Use Pruned Tree, Use Decisions No, Use Priors No, Exhaustive 5000, Node Sample 20000, Bootstrap 0, Method Assessment, Number of Leaves 10, Assessment Mean/Average Square Error 0.000000, Assessment Pratio(0,25).
- Diagram Area:** Results - Node: Probability tree Diagram: Employee. The tree structure is as follows:
 - Root Node (Node Id: 1): industry
 - POWERGENERATION, RETAIL... (Node Id: 2)
 - Node Id: 2: Statistic Train Validation O: 56.81% I: 41.69% Count: 351
 - Node Id: 3: Statistic Train Validation O: 56.74% I: 65.26% Count: 213
 - BANKS, CONSULT, STATE, E... (Node Id: 3)
 - Node Id: 3: Statistic Train Validation O: 56.74% I: 65.26% Count: 224
 - POWERGENERATION, RETAIL... (Node Id: 2)
 - Node Id: 2: Statistic Train Validation O: 56.81% I: 41.69% Count: 351
 - Node Id: 3: Statistic Train Validation O: 56.74% I: 65.26% Count: 213
 - BANKS, CONSULT, STATE, E... (Node Id: 3)
 - Node Id: 3: Statistic Train Validation O: 56.74% I: 65.26% Count: 224
 - profession
 - HR, ETC, BUSINESSDEVELO... (Node Id: 4)
 - Node Id: 4: Statistic Train Validation O: 60.00% I: 37.16% Count: 286
 - Node Id: 5: Statistic Train Validation O: 59.84% I: 65.45% Count: 55
 - COMMERCIAL, SALES, FINAN... (Node Id: 5)
 - Node Id: 5: Statistic Train Validation O: 59.84% I: 65.45% Count: 61
 - profession
 - HR, ETC, BUSINESSDEVELO... (Node Id: 4)
 - Node Id: 4: Statistic Train Validation O: 60.00% I: 37.16% Count: 286
 - Node Id: 5: Statistic Train Validation O: 59.84% I: 65.45% Count: 55
 - COMMERCIAL, SALES, FINAN... (Node Id: 5)
 - Node Id: 5: Statistic Train Validation O: 59.84% I: 65.45% Count: 61
 - way
 - FOOT (Node Id: 6)
 - Node Id: 6: Statistic Train Validation O: 54.88% I: 45.16% Count: 31
 - Node Id: 7: Statistic Train Validation O: 57.11% I: 68.28% Count: 182
 - BUS, CAR Or Missing (Node Id: 7)
 - Node Id: 7: Statistic Train Validation O: 57.11% I: 68.28% Count: 194
 - way
 - FOOT (Node Id: 6)
 - Node Id: 6: Statistic Train Validation O: 54.88% I: 45.16% Count: 31
 - Node Id: 7: Statistic Train Validation O: 57.11% I: 68.28% Count: 182
 - BUS, CAR Or Missing (Node Id: 7)
 - Node Id: 7: Statistic Train Validation O: 57.11% I: 68.28% Count: 194
 - stag
 - < 4.6817 (Node Id: 8)
 - Node Id: 8: Statistic Train Validation O: 100.00% I: 0.00% Count: 8
 - >= 4.6817 Or Missing (Node Id: 9)
 - Node Id: 9: Statistic Train Validation O: 44.00% I: 66.00% Count: 24
 - stag
 - < 4.6817 (Node Id: 8)
 - Node Id: 8: Statistic Train Validation O: 100.00% I: 0.00% Count: 8
 - >= 4.6817 Or Missing (Node Id: 9)
 - Node Id: 9: Statistic Train Validation O: 44.00% I: 66.00% Count: 24

Enterprise Miner - BA706DEMO

File Edit View Actions Options Window Help

BA706DEMO

- Data Sources
- Diagrams
 - Employee
 - HOME
 - Loan
 - Predictive Analytics
 - Students
 - Model Packages

Properties

Property	Value
Pruned values	Use in search
Use Input Once	No
Maximum Branch	2
Maximum Depth	6
Minimum Categor	5
Node Size	5
Number of Rules	5
Number of Surro	0
Split Size	1
Split Rule	Information
Use Decisions	No
Use Priors	No
Exhaustive	5000
Sample	20000
Method	Assessment
Number of Leaves	{}
Assessment Meas	Average Square E
Assessment Fract	0.25

General

Diagram Employee opened

Type here to search

SAS I.T. 2022... Balan... Enter... Enter... Result... ENG 4:53 PM US 12/12/2022

Enterprise Miner - BA706DEMO

File Edit View Actions Options Window Help

BA706DEMO

- Data Sources
- Diagrams
 - Employee
 - HOME
 - Loan
 - Predictive Analytics
 - Students
 - Model Packages

Properties

Property	Value
Pruned values	Use in search
Use Input Once	No
Maximum Branch	2
Maximum Depth	6
Minimum Categor	5
Node Size	5
Number of Rules	5
Number of Surro	0
Split Size	1
Split Rule	Information
Use Decisions	No
Use Priors	No
Exhaustive	5000
Sample	20000
Method	Assessment
Number of Leaves	{}
Assessment Meas	Average Square E
Assessment Fract	0.25

General

Diagram Employee opened

Type here to search

SAS I.T. 2022... Balan... Enter... Enter... Result... ENG 4:54 PM US 12/12/2022

Results - Node: Probability tree Diagram: Employee

Fit Statistics

Target	Target Label	Fit Statistics	Statistics Label	Train	Validation	Test
event	NOBS		Sum of Frequencies	564	565	.
event	MISC		Misclassification Rate	0.347518	0.385841	.
event	MAX		Maximum Absolute Error	0.686813	0.686813	.
event	SSE		Sum of Squared Errors	253.0793	267.8353	.
event	ASE		Average Squared Error	0.243361	0.237022	.
event	RASE		Root Average Squared...	0.473668	0.485849	.
event	DIV		Divisor for ASE	1128	1130	.
event	DFT		Total Degrees of Free...	564	.	.

Output

```

49 *
50 *-----#
51 * Scope Output
52 *-----#
53 *
54 *-----#
55 *-----#
56 *-----#
57 *-----#
58 *-----#
59 *
60 *
61 Variable Importance
62 *
63 *
64 *-----#
65 *-----#
66 *-----#
67 *-----#
68 *-----#
69 *-----#
70 *-----#
71 *-----#
72 *-----#
73 *
74 Tree Leaf Report
75 *
76 *
77 *-----#
78 *-----#
79 *-----#
80 *-----#
81 *-----#
82 *-----#
83 *-----#
84 *-----#
85 *-----#
86 *
87 *
88 *
89 *
90 *-----#

```

Diagram Employee opened

Type here to search

SAS I.T. 2022... Balan... Enter... Enter... Result... ENG 4:54 PM US 12/12/2022

The tree has five leaves with four variable splits. this reduction is a result of pruning using the average squared error measure. The tree shows that 62.89% of employees in Banks, Consult and State industries who go to work in vehicles are likely to quit their jobs. 56% of employees who have worked greater than 4.7 months in the Bank, Consult, and State industries and walk to the office are likely to quit. 62.3% of Commercial, Sales, and Finance professionals in the power generation, and retail industry are also likely to quit.

The Validation Average Square Error for the two-branch classification Tree is 0.237022. This has a similar validation ASE as the two-branch classification tree model.

THREE BRANCH PROBABILITY TREE

The decision tree node was attached to the data partition node with a maximum branch of three, ‘assessment’ used as the method and ‘average squared error’ as the assessment measure. Below shows the outcome after this node was run with the specifications noted.

Enterprise Miner - BA706DEMO

File Edit View Actions Options Window Help

BA706DEMO

- Data Sources
- Diagrams
- Employee**
- GAME
- Loan
- Predictive Analytics
- Student
- Model Packages

Property Value

- Missing values USE IN SEARCH
- Use Input Once No
- Maximum Branch 3
- Maximum Depth 6
- Minimum Category 5
- Tree Size 5
- Number of Rules 5
- Number of Surrogate 0
- Split Size .
- Split Search
- Use Decisions No
- Use Priors No
- Exhaustive 5000
- Node Sample 20000
- Subtree
- Method Assessment
- Number of Leaves 1
- Assessment Meas Average Square Error
- Assessment Fract 0.25

General

General Properties

Results - Node: Probability tree 3 maxB Diagram: Employee

Score Rankings Overlay: event

Cumulative Lift

Tree

Node Id: 1
Statistic: Train Validation
0: 45.47% 49.39%
1: 50.53% 50.62%
Count: 564 565

industry

Leaf Statistics

Sum

Index

Fit Statistics

Target	Target Label	Fit Statistics	Statistics	Train	Validation	Test
event	NOBS	Sum of Fr...	564	565		
event	MISC	Misclassifi...	0.35461	0.39115		
event	MAX	Maximum ...	0.727273	0.727273		
event	SSE	Sum of S...	256.041	270.7844		
event	ASE	Average ...	0.226987	0.239632		
event	RASE	Root Aver...	0.476431	0.489522		
event	DIV	Divisor for...	1128	1130		
event	DFT	Total Deda...	564			

Treemap

Output

```

1 *...
2 User: 301259629
3 Date: December 10, 2022
4 Time: 01:11:09
5 *...
6 * Training Output
7 *...
8
9
10
11

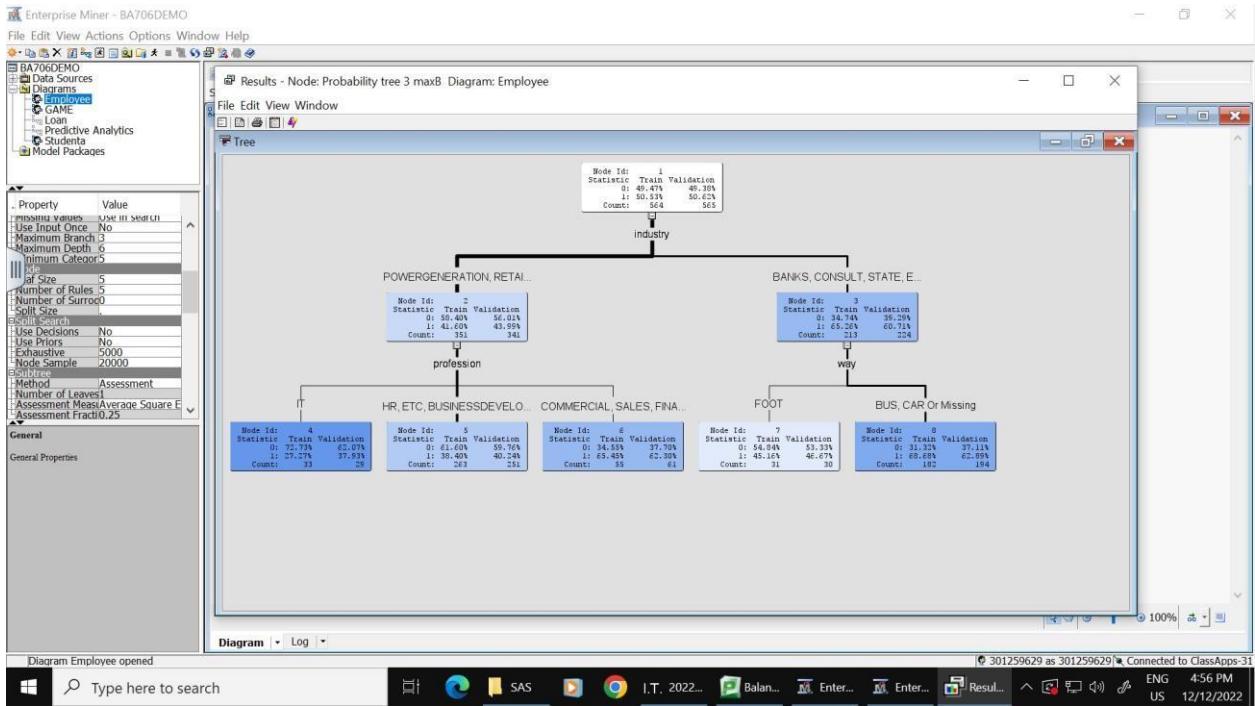
```

Diagram Log

Diagram Employee opened

Type here to search

SAS I.T. 2022... Balan... Enter... Enter... Result... ENG 4:56 PM US 12/12/2022



Enterprise Miner - BA706DEMO

File Edit View Actions Options Window Help

BA706DEMO

Diagrams

- Employee
- IT
- Loan
- Predictive Analytics
- Students
- Model Packages

Property Value

- Input values Use in search
- Use Input Once No
- Maximum Branch 3
- Maximum Depth 6
- Minimum Categor5
- Node Size 5
- Number of Rules 5
- Number of Surro0
- Split Size
- Use Decisions No
- Use Priors No
- Exhaustive 5000
- Sample 20000
- Method Assessment
- Number of Leaves 1
- Assessment Meas/Average Square E
- Assessment Freq0.25

General

General Properties

Results - Node: Probability tree 3 max8 Diagram: Employee

Fit Statistics

Target	Target Label	Fit Statistics	Statistics Label	Train	Validation	Test
event		NBSS	Sum of Frequencies	564	565	
event		MISC	Misclassification Rate	0.35461	0.39115	
event		MAX	Maximum Absolute Error	0.727273	0.727273	
event		SSE	Sum of Squared Errors	256.041	270.7844	
event		ASE	Average Squared Error	0.29307	0.30392	
event		RASE	Root Average Squared...	0.476431	0.489532	
event		DIV	Divisor for ASE	1128	1130	
event		DFT	Total Degrees of Free...	564		

Diagram Employee opened

Type here to search

File Edit View Window

Diagram Employee opened

Type here to search

File Edit View Window

Diagram Employee opened

Type here to search

File Edit View Window

The screenshot shows the SAS Enterprise Miner interface with the following details:

- Enterprise Miner - BA706DEMO** window title.
- Data Sources** pane on the left lists **BA706DEMO**, **Diagrams**, and **Employee**.
- Employee** diagram properties are displayed in the center-left pane, including:
 - Property: **Input values**, Value: **use in score**
 - Use Input Once**: No
 - Maximum Branch**: 3
 - Maximum Node**: 6
 - Minimum Category**: 1
 - Number of Rules**: 3
 - Number of Surrogate**: 0
 - Split Size**: 1
 - Use Decisions**: No
 - Use Priors**: No
 - Exhaustive**: 5000
 - Node Sample**: 20000
 - Method**: Assessment
 - Number of Leaves**: 3
 - Assessment Meas**: Average Square Error
 - Assessment Freq**: 0.25
- Results - Node: Probability tree 3 max8 Diagram: Employee** window title.
- Output** tab selected, displaying the following text output:

```

47 IMT# 1_event Info: event
48
49 -----
50 * Score Output
51 -----
52 -----
53 -----
54 -----
55 * Report Output
56 -----
57 -----
58 -----
59 -----
60 -----
61 Variable Importance
62
63   Number of
64   Splitting
65   Variable  Label  Rules  Importance  Validation
66   Name       Importance  to Training
67   Basis of
68   Industry   1     1.0000  1.0000  1.0000
69   profession 1     0.7410  0.9496  1.1384
70   way        1     0.4443  0.3931  0.8892
71
72 -----
73 -----
74 Tree Leaf Report
75
76   Training    Validation
77   Node Depth Observations Percent  Observations Percent
78
79   1d          1           1
80   5  2       263  0.38    251  0.40
81   8  2       182  0.69    194  0.63
82   6  2       55   0.45    61   0.42
83   4  2       39   0.27    35   0.38
84   7  2       31   0.45    30   0.47
85
86
87
88

```

- Diagram Employee opened** status message.
- Taskbar at the bottom shows various icons and the date/time: **301259629 as 301259629\Connected to ClassApps-31**, **ENG 4:57 PM US 12/12/2022**.

The tree has five leaves with three variable splits. This model has the least amount of significant variables when compared to all the other tree models. The stag variable which shows the amount of time an employee has worked was not considered significant enough to be included. The tree shows that 62.89% of employees in Banks, Consult and State who go to work in vehicles are likely to quit their jobs. 62.3% of Commercial, Sales, and Finance professionals in the power generation industry are likely to quit.

The Validation Average Square Error for the three-branch classification Tree is 0.239632. The Two Branch Probability Tree gives a higher validation assessment than this tree.

EXPLORING OUR DATA

After we were done with analyzing our decision trees we decided to do a state explore to see if we have any missing values in the dataset. As decision trees are forgiving with respect to data, they are not affected by missing values, this step was not necessary before the tree models. However, for the next set of models, regressions and neural networks, the models are unable to run if there are values missing.

The screenshot shows the SAS Enterprise Miner interface with the following details:

- Title Bar:** Enterprise Miner - BA706DEMO
- Menu Bar:** File Edit View Actions Options Window Help
- Toolbars:** Standard toolbar with icons for New, Open, Save, Print, etc.
- Left Panel (Properties):**
 - General tab: Hide Referral Var Yes, Last Run Date 10/12/2022 1:45 AM, Last Run ID fdce004-7610-4x, Last Error, Last Status Complete, Last Run Time 10/12/2022 1:45 AM, Duration 0 Hr, 0 Min, 9.01, Grid Host, User-Added NodeNo.
 - General Properties tab: Create Time 10/12/2022 1:45 AM, Last Run Date 10/12/2022 1:45 AM, Last Run ID fdce004-7610-4x, Last Error, Last Status Complete, Last Run Time 10/12/2022 1:45 AM, Duration 0 Hr, 0 Min, 9.01, Grid Host, User-Added NodeNo.
- Central Window:**
 - Results - Node: StatExplore Diagram: Employee**
 - Interval Variables Table:**

Data Role	Target	Target Level	Variable	Median	Missing	Non Missing	Minimum	Maximum	Mean	Standard Deviation	Skewness	Kurtosis	Role	Label	Scaled Mean Deviation	Maximum Deviation	Level Id
TRAIN	event	0	staq	30.225...	0	279	0.4928...	155.36...	38.761...	35.248...	1.3232...	1.3131...	INPUT	staq	-0.0596...	0.0584...	1
TRAIN	event	1	staq	23.556...	0	285	0.3942...	160.06...	34.442...	32.173...	1.5672...	2.4039...	INPUT	staq	-0.05841...	0.0584...	2
TRAIN	event	0	selfcon...	5.7	0	279	1	10.5787...	1.3180...	0.0032...	-0.72717...	INPUT	selfcon...	0.0302...	0.0295...	1	
TRAIN	event	1	selfcon...	5.7	0	285	1	10.5787...	1.3180...	0.0032...	-0.72717...	INPUT	selfcon...	0.0302...	0.0295...	2	
TRAIN	event	0	independ...	5.5	0	279	1	9.8...	2.2883...	1.7014...	-0.01926...	-0.52105...	INPUT	independ...	-0.02327...	0.0225...	1
TRAIN	event	1	independ...	5.5	0	285	1	10.5537...	1.6874...	0.0385...	-0.16023...	INPUT	independ...	0.0227...	0.0227...	2	
TRAIN	event	0	anxiety	5.6	0	279	1.7	10.5793...	1.7066...	0.0385...	-0.38356...	INPUT	anxiety	0.0165...	0.0164...	1	
TRAIN	event	1	anxiety	5.6	0	285	1.7	10.5793...	1.7066...	0.0385...	-0.38356...	INPUT	anxiety	0.0165...	0.0164...	2	
TRAIN	event	0	novator	6	0	279	1	10.5759...	1.9043...	-0.23808...	-0.51157...	INPUT	novator	-0.01571...	0.0153...	1	
TRAIN	event	1	novator	6	0	285	1	10.5941...	1.8340...	-0.24588...	-0.29084...	INPUT	novator	0.0153...	0.0153...	2	
TRAIN	event	0	extrav... spective...	5.4	0	279	1	10.5720...	1.9076...	-0.07598...	-0.17534...	INPUT	extrav... spective...	-0.00474...	0.00464...	1	
TRAIN	event	1	extrav... spective...	5.4	0	285	1	10.5720...	1.9076...	-0.07598...	-0.17534...	INPUT	extrav... spective...	-0.00474...	0.00464...	2	
TRAIN	event	0	ace	30	0	279	19	54.31230...	6.4133...	0.6127...	-0.00539...	INPUT	ace	0.0030...	0.00301...	1	
TRAIN	event	1	ace	30	0	285	18	54.31041...	7.0992...	0.48201...	-0.41289...	INPUT	ace	-0.00301...	0.00301...	2	

Diagram Employee opened

Type here to search

301259629 as 301259629 Connected to ClassApps 31

ENG 5:00 PM
US 12/12/2022

CAP AND FLOOR

We noticed there were no missing values, hence no impute node required in the diagram. However, there was a skewness in the dataset so we resolved it with a cap and floor by connecting a replacement node. The result shows that the variables ‘age’ and ‘stag’ were affected after running the cap and floor. Age was capped at 51.42497 and floored at 10.84524 while stag was capped and floored at 137.8856 and 64.7267 respectively.

The screenshot displays two windows from the SAS Enterprise Miner interface. The top window is titled 'Results - Node: cap and floor Diagram: Employee' and contains a table titled 'Total Replacement Counts'. The bottom window is also titled 'Results - Node: cap and floor Diagram: Employee' and contains a table titled 'Interval Variables'.

Total Replacement Counts

Variable	Label	Role	Train	Validation
age	age	INPUT	3	4
anxiety	anxiety	INPUT	0	0
extraversion	extraversion	INPUT	0	0
independ	independ	INPUT	0	0
novator	novator	INPUT	0	0
selfcontrol	selfcontrol	INPUT	0	0
stag	stag	INPUT	12	11

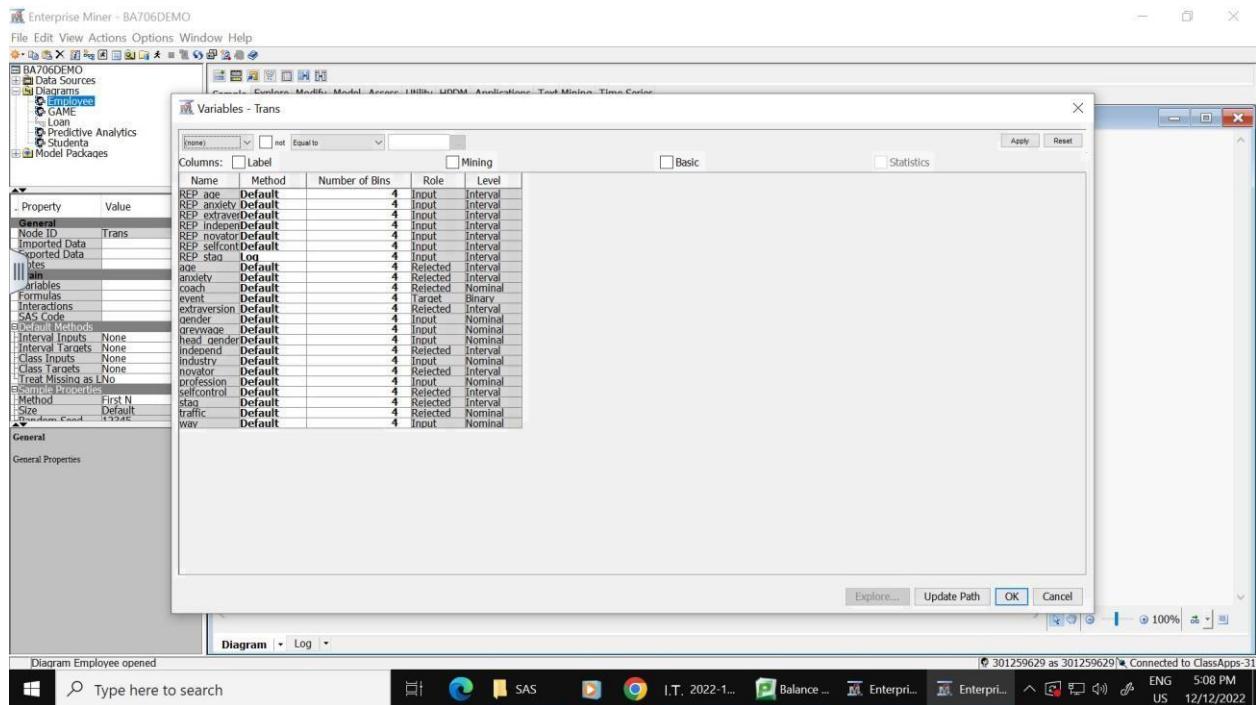
Interval Variables

Variable	Replace Variable	Limits Method	Lower limit	Upper limit	Label	Replacement Method	Lower Replacement Value	Upper Replacement Value
age	REP age	STDDEV	10.84524	51.42497	age	COMPUTED	10.84524	51.42497
anxiety	REP anxiety	STDDEV	0.494046	10.90063	anxiety	COMPUTED	0.494046	10.90063
extraversion	REP extraversion	STDDEV	-0.0169	11.16565	extraversion	COMPUTED	-0.0169	11.16565
independ	REP independ	STDDEV	0.322176	10.50593	independ	COMPUTED	0.322176	10.50593
novator	REP novator	STDDEV	0.242424	11.46041	novator	COMPUTED	0.242424	11.46041
selfcontrol	REP selfcontrol	STDDEV	-0.2212	11.45737	selfcontrol	COMPUTED	-0.2212	11.45737
stag	REP stag	STDDEV	-64.7267	137.8856	stag	COMPUTED	-64.7267	137.8856

A stat explore node was connected to the replacement node to see the effect of the cap and floor on the skewed values.

Data Role	Target	Target Level	Variable	Median	Missing	Non Missing	Minimum	Maximum	Mean	Standard Deviation	Skewness	Kurtosis	Role	Label	Scaled Mean Deviation	Maximum Deviation	Level Id
TRAIN	event	0	REP_staq	30.22587	0	279	0.492813	137.8856	38.53872	34.57223	1.240112	0.943682	INPUT	Replacem...	0.059168	0.057922	1
TRAIN	event	1	REP_staq	23.55647	0	285	0.394251	137.8856	34.27829	31.59778	1.474722	1.907757	INPUT	Replacem...	-0.05792	0.057922	2
TRAIN	event	0	REP_selfcontrol	5.7	0	279	1	10	5.787814	1.978098	0.003512	-0.72717	INPUT	Replacem...	0.030211	0.029575	1
TRAIN	event	1	REP_selfcontrol	5.7	0	285	1	10	5.45193	1.90375	0.132246	-0.32953	INPUT	Replacem...	-0.02958	0.029575	2
TRAIN	event	0	REP_independ	5.5	0	279	1	9.8	5.28853	1.701421	-0.01928	-0.52105	INPUT	Replacem...	-0.02327	0.022782	1
TRAIN	event	1	REP_independ	5.5	0	285	1	10	5.537895	1.687434	0.038587	-0.16023	INPUT	Replacem...	0.022782	0.022782	2
TRAIN	event	0	REP_anxiety	5.6	0	279	1.7	10	5.79319	1.706628	0.087777	-0.38356	INPUT	Replacem...	0.016824	0.016469	1
TRAIN	event	1	REP_anxiety	5.6	0	285	1.7	9.4	5.603509	1.759158	0.217084	-0.71277	INPUT	Replacem...	-0.01647	0.016469	2
TRAIN	event	0	REP_novator	6	0	279	1	10	5.759498	1.904304	-0.23808	-0.51157	INPUT	Replacem...	-0.01571	0.015378	1
TRAIN	event	1	REP_novator	6	0	285	1	10	5.941404	1.834001	-0.24588	-0.29084	INPUT	Replacem...	0.015378	0.015378	2
TRAIN	event	0	REP_extraversion	5.4	0	279	1	10	5.572043	1.907635	-0.07998	-0.17537	INPUT	Replacem...	-0.00474	0.00464	1
TRAIN	event	1	REP_extraversion	5.4	0	285	1	10	5.624561	1.840954	0.059095	-0.52159	INPUT	Replacem...	0.00464	0.00464	2
TRAIN	event	0	REP_aqe	30	0	279	19	51.42497	31.2216	6.382259	0.576385	-0.16957	INPUT	Replacem...	0.003105	0.003039	1
TRAIN	event	1	REP_aqe	30	0	285	18	51.42497	31.03035	7.065542	0.449732	-0.54091	INPUT	Replacem...	-0.00304	0.003039	2

After exploring cap and floor, we realized that the data was still skewed on the Rep stag(which was imputed by the cap and floor to fix the initial skewness of the data). So we used a log transform to try and fix the skewness.



Next was to stat explore the data and we noticed the skewness was sorted as shown below:

Data Role	Target	Target Level	Variable	Median	Missing	Non Missing	Minimum	Maximum	Mean	Standard Deviation	Skewness	Kurtosis	Role	Label	Scaled Mean Deviation	Maximum Deviation	Level Id
TRAIN	event	0	REP_selfcontrol	5.7	0	279	1	10	5.787814	1.978098	0.003512	-0.727171INPUT	Replace...	0.030211	0.29575	1	
TRAIN	event	0	REP_selfcontrol	5.7	0	285	1	10	45193	90375	0.13246	-0.32953INPUT	Replace...	-0.02958	0.29575	2	
TRAIN	event	0	REP_independ	5.5	0	279	1	9	5.62023	1.78728	0.003512	-0.727171INPUT	Replace...	0.022782	0.22782	1	
TRAIN	event	1	REP_independ	5.5	0	285	1	10	5.537895	1.687434	0.038587	-0.16023INPUT	Replace...	0.022782	0.22782	2	
TRAIN	event	0	REP_anxiety	5.6	0	279	1.7	10	5.793319	1.706628	0.087777	-0.38356INPUT	Replace...	0.016824	0.16469	1	
TRAIN	event	0	REP_anxiety	5.6	0	285	1.7	9	5.603509	1.981361	0.217608	-0.12126INPUT	Replace...	-0.016824	0.16469	2	
TRAIN	event	0	REP_novator	6	0	279	1	10	5.933044	1.983044	0.055708	-0.55708	Replace...	0.015167	0.15167	1	
TRAIN	event	1	REP_novator	6	0	285	1	10	5.941404	1.834001	-0.24588	-0.29084INPUT	Replace...	0.015378	0.15378	2	
TRAIN	event	0	LOG REP_stag	3.441247	0	279	0.400662	4.933651	3.220935	0.068798	-0.54579	-0.40702INPUT	Transform...	0.012853	0.012582	1	
TRAIN	event	1	LOG REP_stag	3.200975	0	285	0.332357	4.933651	3.149005	0.068798	-0.54579	-0.40702INPUT	Transform...	-0.012580	0.012582	2	
TRAIN	event	0	REP_extraversion	5.4	0	279	1	10	5.149043	0.076385	-0.07998	-0.17859INPUT	Replace...	0.00464	0.00464	1	
TRAIN	event	1	REP_extraversion	5.4	0	285	1	10	5.624561	1.840954	0.059095	-0.52159INPUT	Replace...	0.00464	0.00464	2	
TRAIN	event	0	REP_aqe	30	0	279	19	51.42497	31.2216	6.382259	0.576385	-0.16957INPUT	Replace...	0.003105	0.003039	1	
TRAIN	event	1	REP_aqe	30	0	285	18	51.42497	31.03035	7.065542	0.449732	-0.54091INPUT	Replace...	-0.00304	0.003039	2	

REPLACEMENT

In logistic regression models, encoding all of the independent variables as dummy variables allows easy interpretation and increases the stability and significance of the coefficients.

Therefore we made the Industry variables **HoReCa** and **etc** ‘unknown’ because we do not know what they are.

The screenshot shows the SAS Enterprise Miner interface with the 'Replacement Editor' dialog box open. The dialog box displays a table of variable replacements. The columns are: Variable, Formatted Value, Replacement Value, Frequency Count, Type, Character Unformatted Value, and Numeric Value. The table includes rows for gender ('head_gender') and industry ('industry'). The 'industry' row shows various categories like Retail, manufacture, IT, Banks, etc., with their corresponding unformatted values and frequencies. The 'etc' row is highlighted with 'Unknown' as the replacement value. The 'HoReCa' row is also highlighted with 'Unknown' as the replacement value. The 'profession' row shows 'HR' as the formatted value and '377C' as the frequency count. The SAS desktop environment is visible in the background, showing other windows and the taskbar.

Variable	Formatted Value	Replacement Value	Frequency Count	Type	Character Unformatted Value	Numeric Value
head_gender	m		299C	m	.	.
head_gender	f		265C	f	.	.
head_gender	_UNKNOWN_	DEFAULT_	-	C	.	.
Industry	Retail		143C	Retail	.	.
Industry	manufacture		76C	manufacture	.	.
Industry	IT		63C	IT	.	.
Industry	Banks		57C	Banks	.	.
Industry	etc	Unknown	49C	etc	.	.
Industry	Consult		37C	Consult	.	.
Industry	State		23C	State	.	.
Industry	PowerGeneration		19C	PowerGeneration	.	.
Industry	transport		19C	transport	.	.
Industry	Building		17C	Building	.	.
Industry	Telecom		16C	Telecom	.	.
Industry	Mining		14C	Mining	.	.
Industry	Pharma		10C	Pharma	.	.
Industry	RealEstate		8C	RealEstate	.	.
Industry	HoReCa	Unknown	7C	HoReCa	.	.
Industry	Agriculture		6C	Agriculture	.	.
Industry	_UNKNOWN_	DEFAULT_	-	C	.	.
profession	HR		377C	HR	.	.

And we combined all the industry variables except for HR, IT and Sales and named it ‘Other’. We also combined the Profession variables BUS and Car and named it ‘Motor’.

Enterprise Miner - BA706DEMO

File Edit View Actions Options Window Help

BA706DEMO

- Data Sources
- Diagrams
- Employee**
 - GAME
 - Loan
 - Predictive Analytics
 - Student
- Model Packages

Sample

Properties

Property	Value
Notes	rain
Interval Variables	
Replacement Edit	
Default Limits MetNone	
Cutoff Values	
Class Variables	
Replacement Edit	
Unknown Levels Tolerance	
Scope	
Replacement Value Computed	
Hide	No
Report	
Replacement RepYes	
Status	
Create Time	14/12/22 3:18 PM
Run ID	7121db5c-673d-4
Last Error	
Last Status	Complete

Class Variables

Replacement options for class variables.

Diag

OK Cancel

Type here to search

SAS I.T. 2022-12-0... Balance for... Enterprise ...

Lesson 8.R Data.csv Show all X

9:07 PM 12/14/2022

9:07 PM 12/14/2022

File Edit View Actions Options Window Help

pod1.centennialcollege.ca/portal/webclient/#/desktop

Replacement Editor-WORK.OUTCLASS

Variable	Formatted Value	Replacement Value	Frequency Count	Type	Character Unformatted Value	Numeric Value
industry	_UNKNOWN_	DEFAULT_	C		.	.
profession	HR		377C	HR	.	.
profession	IT		40C	IT	.	.
profession	Sales		31C	Sales	.	.
profession	etc	Other	22C	etc	.	.
profession	Commercial	Other	14C	Commercial	.	.
profession	Marketing	Other	14C	Marketing	.	.
profession	Consult	Other	13C	Consult	.	.
profession	BusinessDevelopment	Other	12C	BusinessDevelopment	.	.
profession	Finan�e	Other	8C	Finan�e	.	.
profession	Teaching	Other	8C	Teaching	.	.
profession	manage	Other	8C	manage	.	.
profession	Engineer	Other	7C	Engineer	.	.
profession	Accounting	Other	4C	Accounting	.	.
profession	Law	Other	4C	Law	.	.
profession	PR	Other	2C	PR	.	.
profession	_UNKNOWN_	DEFAULT_	C		.	.
way	bus	Motor	341C	bus	.	.
way	car	Motor	157C	car	.	.
way	foot		66C	foot	.	.
way	_UNKNOWN_	DEFAULT_	C		.	.

Model Comparison

Results - Node: Replacement Diagram: Employee

Variable	Formatted Value	Type	Character	Unformatted Value	Binaryic Value	Replacement Value	Label
industry	c	c	etc	.	.	Unknown	
industry	c	c	IndCs	.	.	Unknown	
profession	c	c	etc	.	.	Other	
profession	c	c	Commercial	.	.	Other	
profession	c	c	Marketing	.	.	Other	
profession	c	c	Commerical	.	.	Other	
profession	c	c	BusinessDevelopment	.	.	Other	
profession	c	c	Finanuue	.	.	Other	
profession	c	c	Teaching	.	.	Other	
profession	c	c	Manage	.	.	Other	
profession	c	c	Enginier	.	.	Other	
profession	c	c	Accounting	.	.	Other	
profession	c	c	Law	.	.	Other	
profession	c	c	PR	.	.	Other	
way	c	c	bus	.	.	Motor	
way	c	c	car	.	.	Motor	

Total Replacement Counts

Variable	Role	Label	Train	Validation
industry	INPUT		56	49
profession	INPUT		116	116
way	INPUT		498	514

Dummy variables, which are useful in regression analysis, use the values 0 or 1 to represent the absence or existence of any categorical effect that would be anticipated to shift the outcome. So, now we can begin building our regression models.

REGRESSION

As stated earlier, there were no missing values so we did not need to connect the impute node.

Therefore, we are going into regression. We first created :

FULL REGRESSION

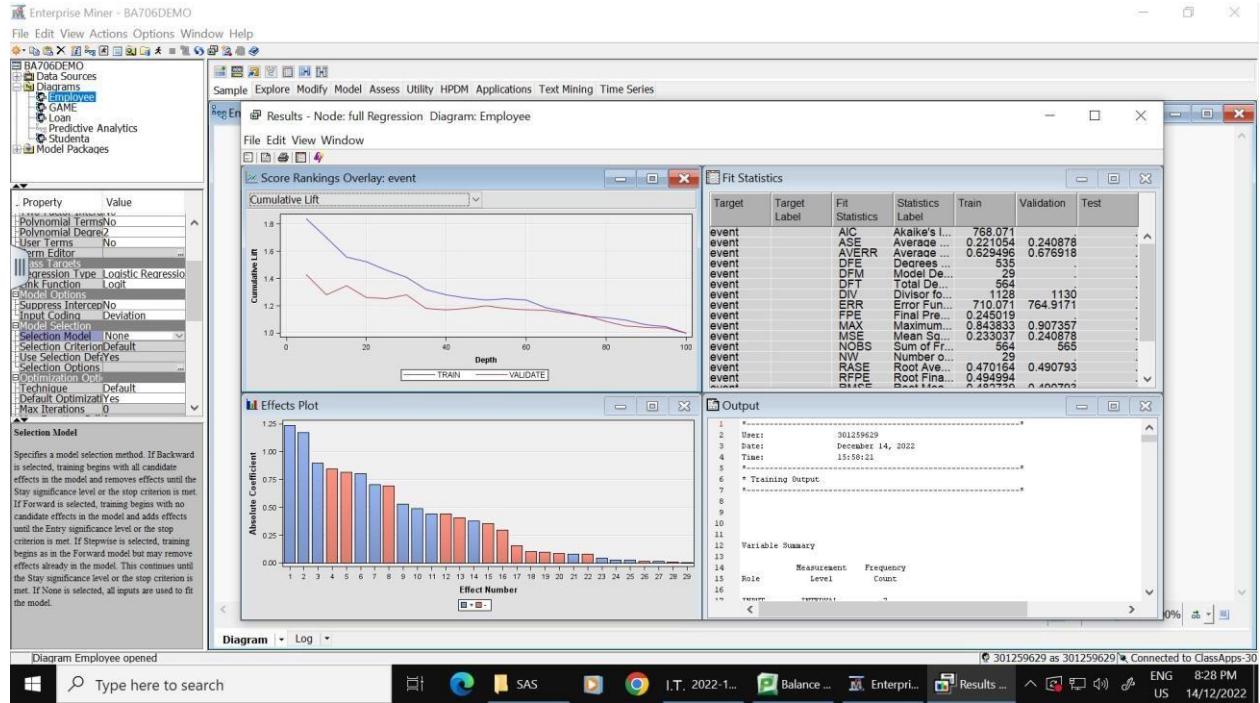
And in other to optimize complexity, we also created:

FORWARD REGRESSION STEPWISE REGRESSION

BACKWARD REGRESSION

FULL REGRESSION

The full regression used 28 variables to run the model. The average squared error is 0.240878



Enterprise Miner - BA706DEMO

File Edit View Actions Options Window Help

BA706DEMO

Data Sources

Diagrams

- Employee
- HOME
- Loans
- Predictive Analytics
- Students
- Model Packages

Property Value

- Polynomial Terms No
- Polynomial Degree 2
- User Terms No
- Term Editor
- Regression Type Logistic Regression
- Link Function Logit
- Model Options
- Output Deviance
- Input Coding Deviation
- Model Selection
- Selection None
- Selection Criterion Gault
- Use Selection Def Yes
- Selection Options ...
- Optimization Ord
- Termination Default
- Default Optimizaties
- Max Iterations 0

Selection Model

Specifies a model selection method. If Backward is selected, training begins with all variables.

Diagram Employee opened

Type here to search

Lesson 8.R Data.csv

Results Node: full Regression Diagram: Employee

Fit Statistics

Target	Target Label	Fit Statistics	Statistics Label	Train	Validation	Test
event		AIC	Akaike's Information Cr...	768.071		
event		ASE	Average Squared Error	0.221054	0.240878	
event		AVGRR	Average Residual Sum of Squares	0.629496	0.676918	
event		DFE	Degrees of Freedom f...	535		
event		DFM	Model Degrees of Free...	29		
event		DFT	Total Degrees of Free...	564		
event		DV	Degrees of Variation	1128	1130	
event		ERR	Error Function	710.071	764.9171	
event		FPE	Final Prediction Error	0.245019		
event		MAX	Maximum Absolute Error	0.843833	0.907357	
event		MSE	Mean Square Error	0.230307	0.240878	
event		NBDS	Sum of Frequencies	664	665	
event		NW	Number of Estimate W...	29		
event		RASE	Root Average Sum of Squares	0.470164	0.490793	
event		RFPE	Root Final Prediction Error	0.484954		
event		RMSE	Root Mean Squared Er...	0.482739	0.490793	
event		SBC	Schwarz's Bayesian Cr...	893.7876		
event		SST	Sum of Squared Errors	249.3494	272.932	
event		SUMW	Sum of Unweighted ...	1128	1130	
event		MISC	Misclassification Rate	0.379433	0.39469	

ted to ClassApps-30

8:31 PM 12/14/2022 ENG US

Enterprise Miner - BA706DEMO

File Edit View Actions Options Window Help

BA706DEMO

Data Sources

Diagrams

- Employee
- HOME
- Loans
- Predictive Analytics
- Students
- Model Packages

Property Value

- Polynomial Terms No
- Polynomial Degree 2
- User Terms No
- Term Editor
- Regression Type Logistic Regression
- Link Function Logit
- Model Options
- Output Deviance
- Input Coding Deviation
- Model Selection
- Selection None
- Selection Criterion Gault
- Use Selection Def Yes
- Selection Options ...
- Optimization Ord
- Termination Default
- Default Optimizaties
- Max Iterations 0

Selection Model

Specifies a model selection method. If Backward is selected, training begins with all variables.

Diagram Employee opened

Type here to search

Lesson 8.R Data.csv

Results Node: full Regression Diagram: Employee

Output

Effect	Point Estimate
REF_star	0.907
REF_age	1.003
REF_education	0.926
REF_extraversion	0.918
REF_independ	1.012
REF_industry Agriculture vs transport	6.503
REF_industry Banks vs transport	4.113
REF_industry Banks vs transport	9.346
REF_industry Commute vs transport	4.478
REF_industry IT vs transport	1.169
REF_industry Mining vs transport	3.858
REF_industry Pharms vs transport	5.904
REF_industry RealEstate vs transport	1.127
REF_industry RealEstate vs transport	1.136
REF_industry Retail vs transport	1.857
REF_industry States vs transport	2.721
REF_industry Telecom vs transport	2.279
REF_industry Unknown vs transport	2.712
REF_industry Manufacturing vs transport	1.764
REF_motorcar	1.045
REF_profession HR vs Sales	0.749
REF_profession IT vs Sales	0.671
REF_profession Other vs Sales	2.039
REF_selfcontrol	0.905
REF_way Motor vs foot	2.667
gender F vs u	1.171
greyage grey vs white	0.572
head_pendet F vs u	0.995

ted to ClassApps-30

8:32 PM 12/14/2022 ENG US

Interpretation of Full Regression

Variables		point Estimate	When compared with the Transport industry
-----------	--	----------------	---

REP_industry	Building vs transport	8.546	Employees in the Building industry are 8.5 times more likely to quit.
REP_industry	Agric vs transport	6.503	Employees in the Agriculture Industry are 6.5 times more likely to quit.
REP_industry	pharma vs transport	5.904	Employees in the Pharmacy industry are 5.9 times more likely to quit.
REP_industry	Consult vs transport	4.478	Employees in the Consult industry are 4.5 times more likely to quit.
REP_industry	Banks vs transport	4.113	Employees in the Bank industry are 4.1 times more likely to quit.
REP_industry	Mining vs transport	3.858	Employees in the Mining industry are 3.9 times more likely to quit.
REP_industry	State vs transport	2.721	Employees in the State industry are 2.7 times more likely to quit.

REP_industry	Unknown vs transport	2.712	Employees in the Unknown industry are 2.7 times more likely to quit.
REP_industry	Telecom vs transport	2.270	Employees in the Telecom industry are 2.3 times more likely to quit.
REP_industry	Retail vs transport	1.857	Employees in the Retail industry are 85.7% more likely to quit.
REP_industry	Manufacture vs transport	1.704	Employees in the Manufacture industry are 70.4% more likely to quit.
REP_industry	PowerGen vs transport	1.327	Employees in the PowerGen industry are 32.7% more likely to quit.

REP_industry	IT vs transport	1.169	Employees in the IT industry are 16.9% more likely to quit.
REP_industry	RealEstate vs transport	1.136	Employees in the Real Estate industry are 13.6% more likely to quit.
			When compared with Sales Profession
REP_profession	Other vs Sales	2.039	Employees in other professions besides HR and IT are more likely to quit
REP_profession	HR vs Sales	0.749	Employees in HR are 25.1% less likely to quit.

REP_profession	IT vs Sales	0.671	Employees in IT are 32.9% less likely to quit.
Other Variables			

REP_way	Motor vs foot	2.667	Employees who come by Motor are 2.7 times more likely quit than those who come by foot
gender	f vs m	1.171	Female employees are 17.1% more likely to quit than male employees

REP_novator		1.045	For each 1 score the odds of quitting changed by a factor 1.045, a 4.5% increase
dependent		1.012	For each 1 score the odds of quitting changed by a factor of 1.012, a 1.2% increase

REP_age		1.003	For every year added in age, employees are 0.3% more quit
head_gender	f vs m	0.985	Female supervisor is 2% less likely to quit than a male supervisor
greywage	grey vs white	0.972	Employees who earn above minimum wage are 2.8% less likely to quit than those who earn minimum wage
REP_anxiety		0.926	For each 1 score the odds of quitting changed by a factor of 0.926, a 7.4% decrease
raversion		0.918	For each 1 score the odds of quitting changed by a factor of 0.918, a 8.2% decrease
LOG_REP_stag		0.907	For each 1 score the odds of quitting changed by a factor of 0.907, a 9.3% decrease
REP_selfcontrol		0.905	For each 1 score the odds of quitting changed by a factor of 0.905, a 9.5% decrease

FORWARD REGRESSION

The forward regression used 18 variables to build the model and the outcome is an ASE of 0.237256 which is a better result than the full regression.

The gender of the employee or supervisor, self control score , anxiety score, all other scores, wage level, and experience were not considered significant in building this model.

Enterprise Miner - BA706DEMO

File Edit View Actions Options Window Help

Score Rankings Overlay: event

Cumulative Lift

Fit Statistics

Target	Target Label	Fit Statistics	Statistics Label	Train	Validation	Test
event		AIC	Akaike's I...	756.7345		
event		ASE	Average ...	0.224052	0.237256	
event		AVERR	Average ...	0.637176	0.667418	
event		DFE	Degrees...	545		
event		DFM	Model De...	19		
event		DFT	Total De...	564		
event		DIV	Divisor fo...	1128	1130	
event		ERR	Error Fun...	718.7345	754.1829	
event		FPE	Final Pre...	0.239674		
event		MAX	Maximum...	0.861805	0.877769	
event		MSE	Mean Sq...	0.231863	0.237256	
event		NOBS	Sum of Fr...	564	565	
event		NW	Number o...	19		
event		RASE	Root Ave...	0.473342	0.48709	

Effects Plot

Output

```

1 *-----*
2 User: 301259629
3 Date: December 14, 2022
4 Time: 15:59:23
5 *-----*
6 * Training Output *
7 *-----*
8
9
10
11
12 Variable Summary
13
14 Measurement Frequency
15 Role Level Count

```

Type here to search

I.T. 2022-1... SAS Chrome I.T. 2022-1... Enterprise Miner Results... Show all

ENG 8:36 PM
US 14/12/2022

8:36 PM
12/14/2022

Enterprise Miner - BA706DEMO

File Edit View Actions Options Window Help

BA706DEMO

- Data Sources
- Diagrams
 - Employee**
 - HOME
 - Loans
 - Predictive Analytics
 - Students
- Model Packages

Property Value

Polynomial Terms No
Polynomial Degree
User Terms No
Term Editor
Regression Type Logistic Regression
Link Function Logit
Model Options
Input Coding Deviation
Model Selection

- Selection Model Forward
- Use Selection Error Validation Error
- Use Selection Def Yes
- Selection Options ...

- Optimization Options
- Termination Rule Default
- Default Optimizatives
- Max Iterations 0
- General
- General Properties

Diagram Employee opened

Type here to search

SAS I.T. 2022-1... Balance... Enterprise... Results... ENG 8:36 PM 12/14/2022

Lesson 8.R Data.csv Show all

0% ↻

Results - Node: forward reg Diagram: Employee

Fit Statistics

Target	Target Label	Fit Statistics	Statistics Label	Train	Validation	Test
event		AIC	Akaike's Information Cr...	756.7345		
event		ASE	Average Squared Error	0.224052	0.237256	
event		AVERR	Average Error per iteration	0.637176	0.667418	
event		DFE	Degrees of Freedom f...	545		
event		DFM	Model Degrees of Free...	19		
event		DFT	Total Degrees of Free...	564		
event		DV	Degrees of Freedom	1128	1130	
event		ERR	Error Function	718.7345	754.1829	
event		FPE	Final Prediction Error	0.239674		
event		MAX	Maximum Absolute Error	0.861805	0.877769	
event		MSE	Mean Square Error	0.231863	0.237256	
event		NBDS	Sum of Frequencies	564	565	
event		NW	Number of Estimate W...	19		
event		RASE	Root Average Sum of ...	0.473342	0.48709	
event		RFPE	Root Final Prediction	0.488955		
event		RMSE	Root Mean Squared E...	0.481522	0.48709	
event		SBC	Schwarz's Bayesian Cr...	839.1005		
event		SST	Sum of Squared Errors	252.731	268.0986	
event		SUMW	Sum of Weights ...	1128	1130	
event		MISC	Misclassification Rate	0.37234	0.410619	

0% ↻

Enterprise Miner - BA706DEMO

File Edit View Actions Options Window Help

BA706DEMO

- Data Sources
- Diagrams
 - Employee**
 - HOME
 - Loans
 - Predictive Analytics
 - Students
- Model Packages

Property Value

Polynomial Terms No
Polynomial Degree
User Terms No
Term Editor
Regression Type Logistic Regression
Link Function Logit
Model Options
Input Coding Deviation
Model Selection

- Selection Model Forward
- Use Selection Error Validation Error
- Use Selection Def Yes
- Selection Options ...

- Optimization Options
- Termination Rule Default
- Default Optimizatives
- Max Iterations 0
- General
- General Properties

Diagram Employee opened

Type here to search

SAS I.T. 2022-1... Balance... Enterprise... Results... ENG 8:36 PM 12/14/2022

Lesson 8.R Data.csv Show all

0% ↻

Results - Node: forward reg Diagram: Employee

Output

Sample	REF_profession	Other	1	0.45325	0.2176	8.99	0.0027	1.920
444	REF_way	Motor	1	0.46483	0.1496	9.63	0.0019	1.591
445								
446								
447								
448								
449								
450								
451	Effect							
452								
453	REF_industry	Agriculture vs transport						
454	REF_industry	Banks vs transport						
455	REF_industry	Building vs transport						
456	REF_industry	Consult vs transport						
457	REF_industry	IT vs transport						
458	REF_industry	Manufacture vs transport						
459	REF_industry	Pharm vs transport						
460	REF_industry	PowerGeneration vs transport						
461	REF_industry	RealEstate vs transport						
462	REF_industry	Retail vs transport						
463	REF_industry	Services vs transport						
464	REF_industry	Telecom vs transport						
465	REF_industry	Wholesom vs transport						
466	REF_profession	Manufacture vs transport						
467	REF_profession	IT vs Sales						
468	REF_profession	IT vs Services						
469	REF_profession	Other vs Sales						
470	REF_profession	Other vs Services						
471	REF_way	Motor vs foot						
472								
473								
474								
475								
476								
477								
478								
479								
480								
481								
482								
483								
484								
485								
486								
487								
488								
489								
490								
491								
492								
493								
494	*	Scope Output	*					
495	*	-----	*					
496	*	-----	*					
497	*	-----	*					
498	*	Report Output	*					

0% ↻

Interpretation of Forward Regression

		PoinT Estimate	When compared with the Transport industry
REP_industry	Building vs transport	7.635	Employees in the Building industry are 7.6 times more likely to quit.
REP_industry	Agric vs transport	6.769	Employees in the Agriculture industry are 6.8 times more likely to quit.
REP_industry	Pharma vs transport	5.034	Employees in the Pharmaceutical industry are 5 times more likely to quit.
REP_industry	Banks vs transport	4.600	Employees in the Banking industry are 4.6 times more likely to quit.

REP_industry	Consult vs transport	4.408	Employees in the Consulting industry are 4.4 times more likely to quit.
REP_industry	Mining vs transport	3.637	Employees in the Mining industry are 3.6 times more likely to quit.
REP_industry	State vs transport	3.063	Employees in the State industry are 3 times more likely to quit.

REP_industry	Unknown vs transport	2.606	Employees in the all other industries not listed are 2.6 times more likely to quit.
REP_industry	Telecom vs transport	2.356	Employees in the Telecom industry are 2.4% more likely to quit.
REP_industry	Retail vs transport	1.723	Employees in the Retail industry are 72.3% more likely to quit.
			Employees in the Manufacture industry are 62.8% more likely to quit.

REP_industry	Manufacture vs transport	1.628	
REP_industry	PowerGen vs sport	1.292	Employees in the Power generating industry are 29.2% times more likely to quit.
REP_industry	RealEstate vs transport	1.291	Employees in the Real Estate industry are 29.1% times more likely to quit.
REP_industry	IT vs transport	0.880	Employees in the IT industry are 12% less likely to quit.
			When compared with Sales Profession
REP_profession	HR vs Sales	1.762	Employees in professions other than HR and IT are 76% more likely to quit
REP_profession	HR vs Sales	0.715	Employees in HR are 28.5% less likely to quit
REP_profession	IT vs Sales	0.563	Employees in IT are 43.7% less likely to quit
			Other Variables

REP_way	Motor vs foot	2.531	loyees who go to the office with a vehicle are 2.5 times more likely to quit than those who walk to the office.
---------	------------------	-------	---

STEPWISE REGRESSION

The stepwise model also used 18 variables and had the same outcome as the forward regression with an ASE of 0.237256.

Enterprise Miner - BA706DEMO

File Edit View Actions Options Window Help

Score Rankings Overlay: event

Cumulative Lift

Target	Target Label	Fit Statistics	Statistics Label	Train	Validation	Test
event	AIC	Akaike's I...	756.7345			
event	ASE	Average ...	0.224052	0.237256		
event	AVERR	Average ...	0.637176	0.667418		
event	DFE	Degrees ...	545			
event	DFM	Model De...	19			
event	DFT	Total De...	564			
event	DIV	Divisor fo...	1128	1130		
event	ERR	Error Fun...	718.7345	754.1829		
event	FPE	Final Pre...	0.239674			
event	MAX	Maximum...	0.861805	0.877769		
event	MSE	Mean Sq...	0.231863	0.237256		
event	NOBS	Sum of Fr...	564	565		
event	NW	Number o...	19			
event	RASE	Root Ave...	0.473342	0.48709		

Effects Plot

Output

```

1 -----
2 User: 301259629
3 Date: December 14, 2022
4 Time: 15:33:05
5 -----
6 * Training Output
7 -----
8
9
10
11
12 Variable Summary
13
14 Measurement Frequency
15 Role Level Count

```

Type here to search

Lesson 8.R Data.csv

Show all

8:37 PM 12/14/2022

The screenshot displays a dual-monitor setup. The left monitor shows the Enterprise Miner software interface, specifically the 'Results - Node: Regression (stepwise) Diagram: Employee' window. This window contains a 'Fit Statistics' table and a 'Output' section with 'Oddsr Ratios Estimates'. The right monitor also shows the same Enterprise Miner interface, mirroring the left one. The taskbar at the bottom of both monitors includes icons for Data.csv, SAS, Chrome, I.T. 2022-1..., Balance..., Enterprise..., Results..., and various system icons. The system tray shows the date as 14/12/2022 and the time as 8:38 PM.

Interpretation of Stepwise Regression

Variables		Point Estimate	When compared to Transport Industry
-----------	--	----------------	--

Rep_Industry	Agriculture vs transport	6.769	Employees in the Agriculture industry are 6.8 times more likely to quit
Rep_Industry	Banks vs transport	4.600	Employees in the Bank industry are 4.6 times more likely to quit
Rep_Industry	Building vs transport	7.635	Employees in the Building industry are 7.6 times more likely to quit
Rep_Industry	Consult vs transport	4.408	Employees in the Consult industry are 4.4 times more likely to quit
Rep_Industry	IT vs transport	1.198	Employees in the IT industry are 19.8% more likely to quit
Rep_Industry	Mining vs transport	3.637	Employees in the Mining industry are 3.6 times more likely to quit
Rep_Industry	Pharma vs transport	5.034	Employees in the Pharma industry are 5 times more likely to quit

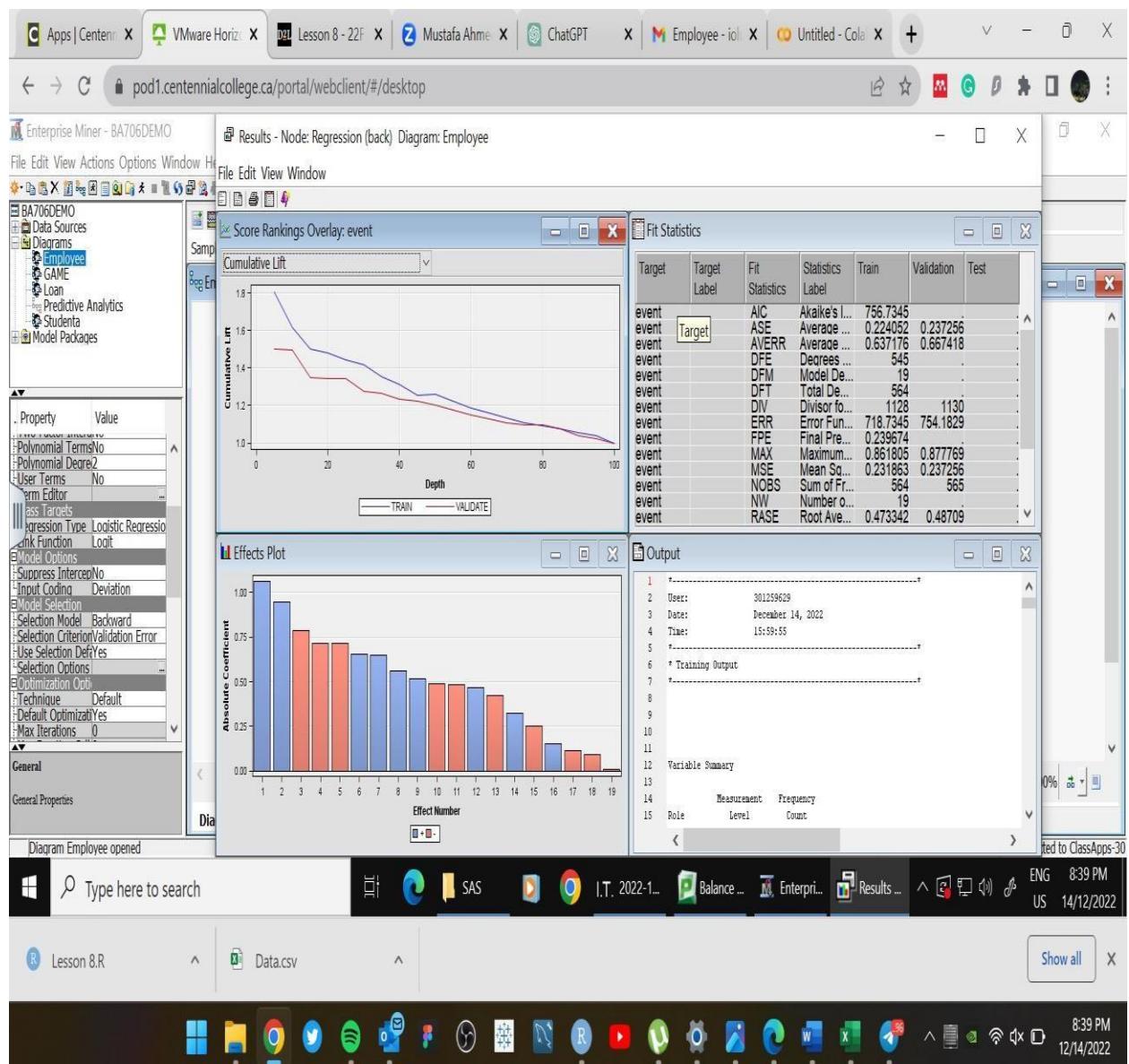
Rep_Industry	Power Generation vs transport	1.292	Employees in the Power Generation industry are 29.2% more likely to quit
Rep_Industry	Real Estate vs transport	1.291	Employees in the Real Estate industry are 29.1% more likely to quit

Rep_Industry	Retail vs transport	1.723	Employees in the Retail Industry are 72.3% more likely to quit
Rep_Industry	State vs transport	3.063	Employees in the State industry are 3.1 times more likely to quit
Rep_Industry	Telecom vs transport	2.356	Employees in the Telecom industry are 2.4 times more likely to quit
Rep_Industry	Unknown vs transport	2.606	Employees in the Unknown industry are 2.6 times more likely to quit
Rep_Industry	Manufacture vs transport	1.628	Employees in the Manufacture industry are 62.8% more likely to quit
			When compared with Sales Profession

Rep_Profession	HR vs Sales	0.715	Employees in HR are 28.5% less likely to quit
Rep_Profession	IT vs Sales	0.563	Employees in IT are 43.7% less likely to quit
Rep_Profession	Other vs Sales	1.762	Employees in Other Professions are 76.2% more likely to quit
			Other Variable
Rep_Way	Motor vs Foot	2.531	Employees who come to work by Motor are 2.5 times more likely to quit than those who come by Foot

BACKWARD REGRESSION

The backward regression model also used 18 variables and had the same outcome as the forward regression and stepwise regression with an ASE of 0.237256.



The screenshot shows a Microsoft Windows desktop environment with the following details:

- Browser Tabs:** C Apps | Centenial, VMware Horiz., Lesson 8 - 22f, Mustafa Ahme, ChatGPT, Employee - iol, Untitled - Cola.
- SAS Application:** Enterprise Miner - BA706DEMO. The interface includes a left sidebar with project navigation (Data Sources, Diagrams, Model Packages), a main pane titled "Fit Statistics" showing regression results, and a bottom pane titled "Odds Ratio Estimates".
- Taskbar:** Shows icons for File Explorer, Microsoft Edge, SAS, Google Chrome, I.T. 2022-1..., Balance..., Enterprise..., Results..., and several system icons like battery, signal, and volume.
- System Status:** ENG 8:40 PM US 14/12/2022.

Interpretation of Backward Regression

Variables		Point Estimate	When compared to Transport Industry
Rep_Industry	Agriculture vs transport	6.769	Employees in the Agriculture industry are 6.8 times more likely to quit
Rep_Industry	Banks vs transport	4.600	Employees in the Bank industry are 4.6 times more likely to quit
Rep_Industry	Building vs transport	7.635	Employees in the Building industry are 7.6 times more likely to quit
Rep_Industry	Consult vs transport	4.408	Employees in the Consult industry are 4.4 times more likely to quit
Rep_Industry	IT vs transport	1.198	Employees in the IT industry are 19.8% more likely to quit
Rep_Industry	Mining vs transport	3.637	Employees in the Mining industry are 3.6 times more likely to quit

Rep_Industry	Pharma vs transport	5.034	Employees in the Pharma industry are 5 times more likely to quit
Rep_Industry	Power Generation vs transport	1.292	Employees in the Power Generation industry are 29.2% more likely to quit

Rep_Industry	Real Estate vs transport	1.291	Employees in the Real Estate industry are 29.1% more likely to quit
Rep_Industry	Retail vs transport	1.723	Employees in the Retail Industry are 72.3% more likely to quit
Rep_Industry	State vs transport	3.063	Employees in the State industry are 3.1 times more likely to quit
Rep_Industry	Telecom vs transport	2.356	Employees in the Telecom industry are 2.4 times more likely to quit
Rep_Industry	Unknown vs transport	2.606	Employees in the Unknown industry are 2.6 times more likely to quit

Rep_Industry	Manufacture vs transport	1.628	Employees in the Manufacture industry are 62.8% more likely to quit
			When compared with Sales Profession
Rep_Profession	HR vs Sales	0.715	Employees in HR are 28.5% less likely to quit
Rep_Profession	IT vs Sales	0.563	Employees in IT are 43.7% less likely to quit
Rep_Profession	Other vs Sales	1.762	Employees in Other Professions are 76.2% more likely to quit
			Other Variable
Rep_Way	Motor vs Foot	2.531	Employees who come to work by Motor are 2.5 times more likely to quit than those who come by Foot

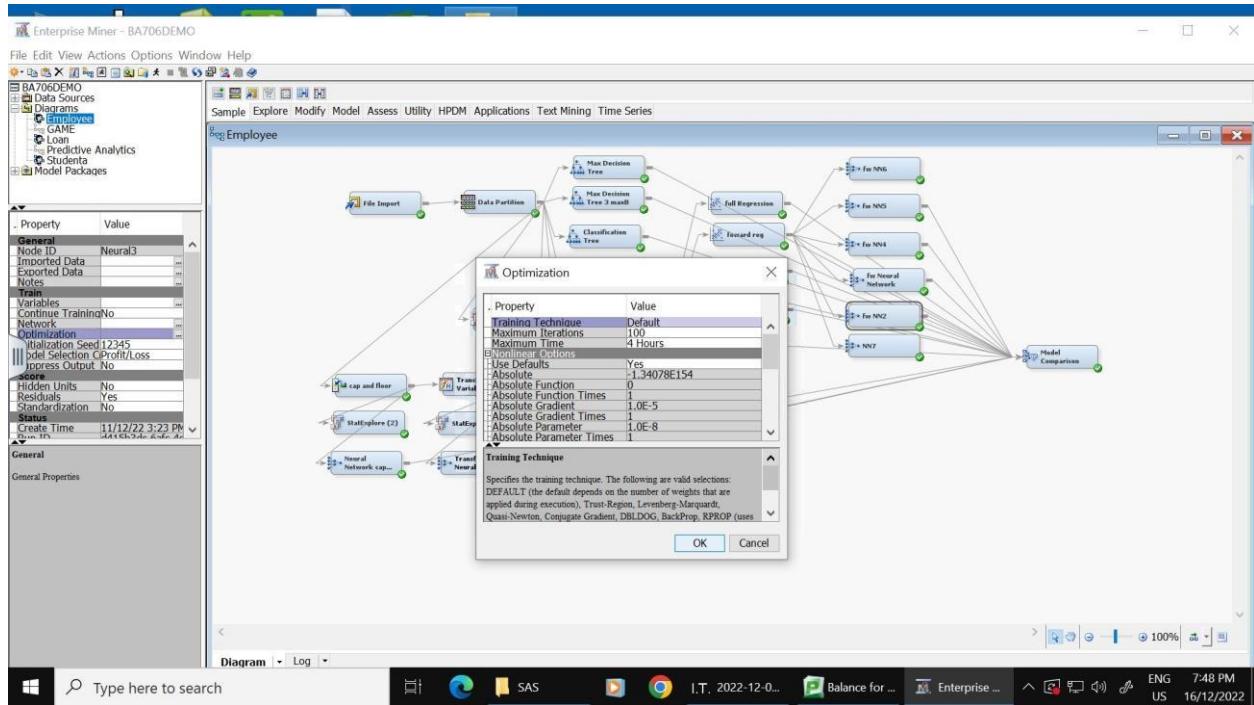
NEURAL NETWORK

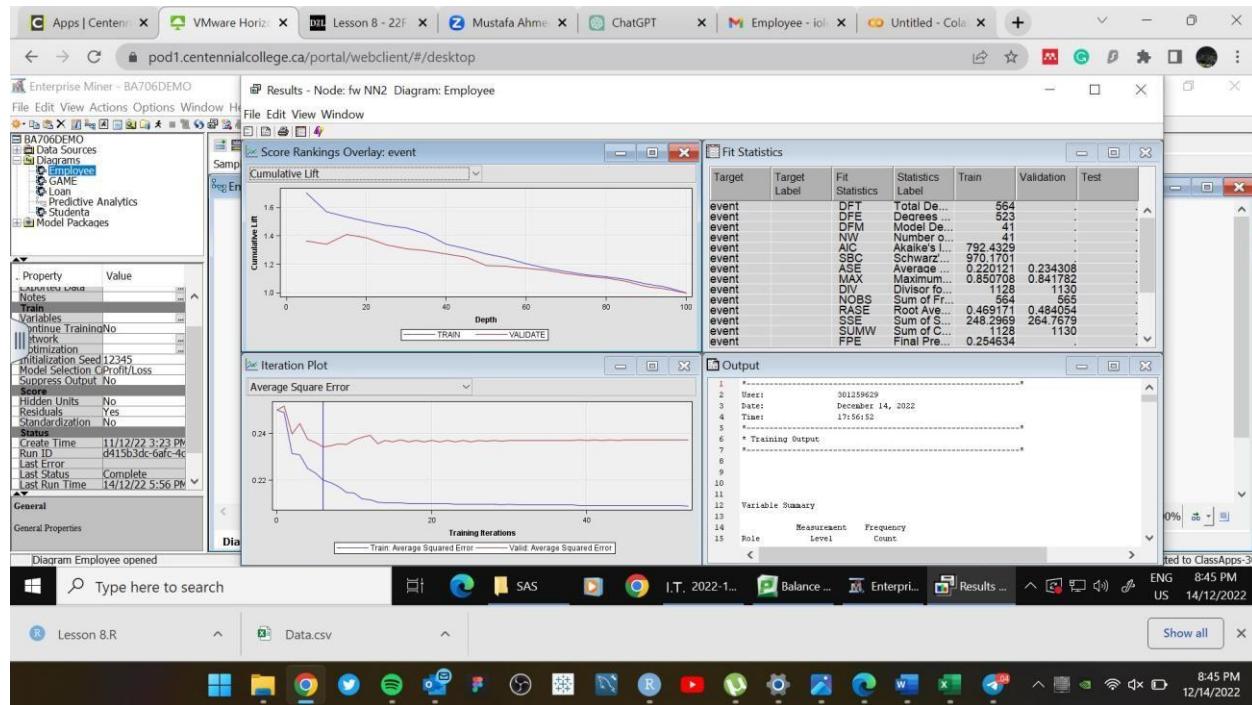
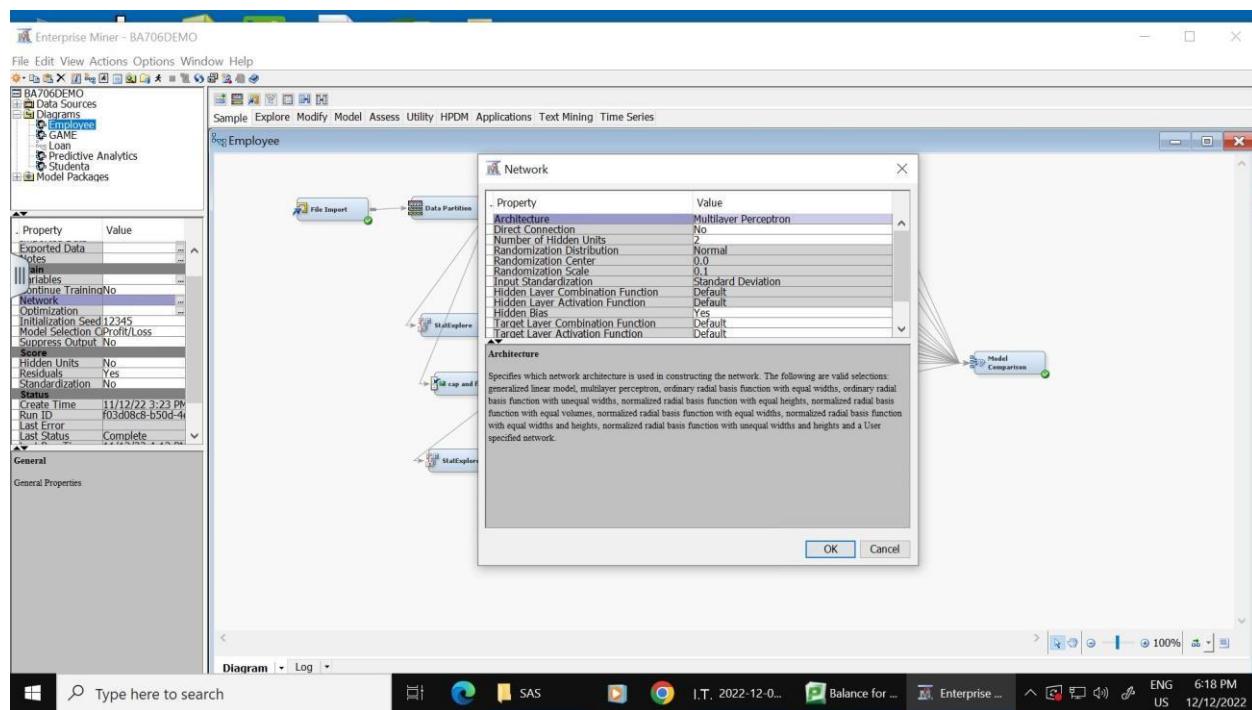
With the full regression having an ASE of 0.240878 and Forward, Stepwise , and Backward regression models all having the same ASE of 0.237256 we have an option of choosing from any of the three with the better average squared error. We chose the Forward regression model because of its additive style of selecting variables.

All neural networks were run using 100 iterations, and preliminary training was turned off. The distinguishing factor for all the neural networks run was the number of hidden units used.

NEURAL NETWORK WITH 2 HIDDEN UNITS

This neural network node was run using two hidden units. The ASE was 0.234308.

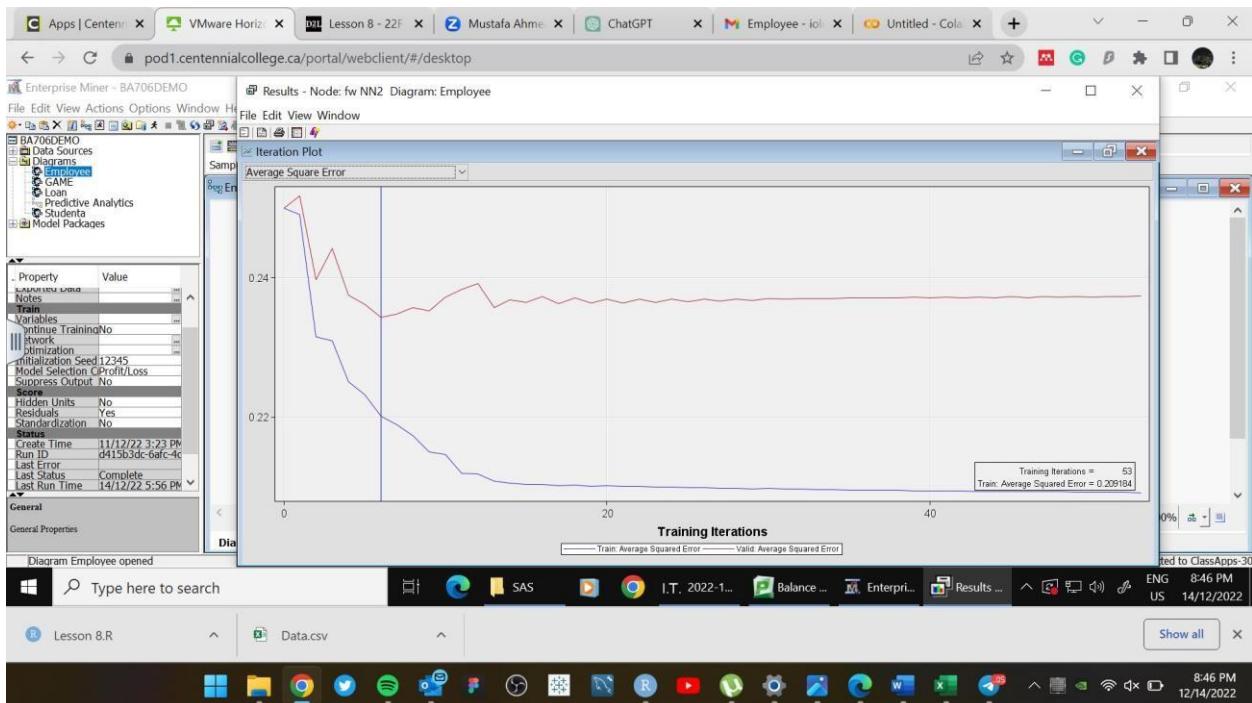




Screenshot of SAS Enterprise Miner showing the 'Fit Statistics' table for a neural network model named 'fw NN2'. The table compares Fit Statistics across Train, Validation, and Test datasets.

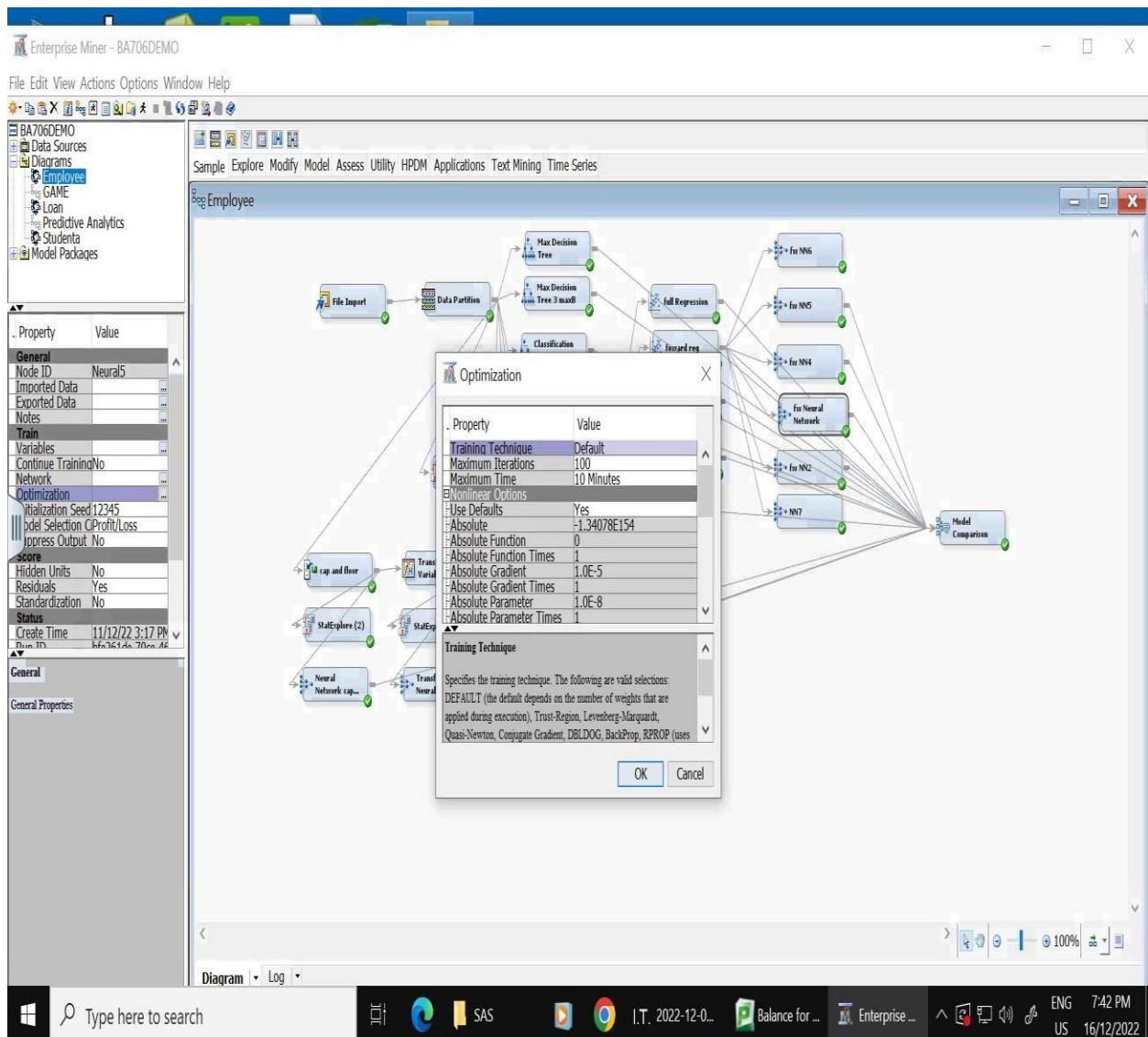
Target	Target Label	Fit Statistics	Statistics Label	Train	Validation	Test
event		DFT	Total Degrees of Free...	564	.	.
event		DFE	Degrees of Freedom f...	523	.	.
event		DNW	Model Degrees of Free...	41	.	.
event		NW	Number of Estimated...	41	.	.
event		AIC	Akaike's Information Cr...	792.4329	.	.
event		BIC	Schwarz's Bayesian Cr...	970.1701	.	.
event		ASE	Average Squared Error	0.234308	0.234308	.
event		MAX	Maximum Absolute Error	0.850708	0.841782	.
event		DIV	Divisor for ASE	1128	1130	.
event		NOBS	Sum of Frequencies	564	565	.
event		RASE	Root Average Squared E...	0.469171	0.484054	.
event		SSE	Sum of Squared Errors	248.2969	264.7679	.
event		SUMW	Sum of Case Weights ...	1128	1130	.
event		FPE	Final Prediction Error	0.234304	0.234308	.
event		MSE	Mean Squared Error	0.234377	0.234308	.
event		RFPE	Root Final Prediction ...	0.504612	0.484054	.
event		RMSE	Root Mean Squared E...	0.487214	0.484054	.
event		AVERR	Average Error Function	0.681518	0.681094	.
event		ERR	Error Function	747.0524	747.0524	.
event		MISC	Misclassification Rate	0.359929	0.39469	.
event		WRONG	Number of Wrong Clas...	203	223	.

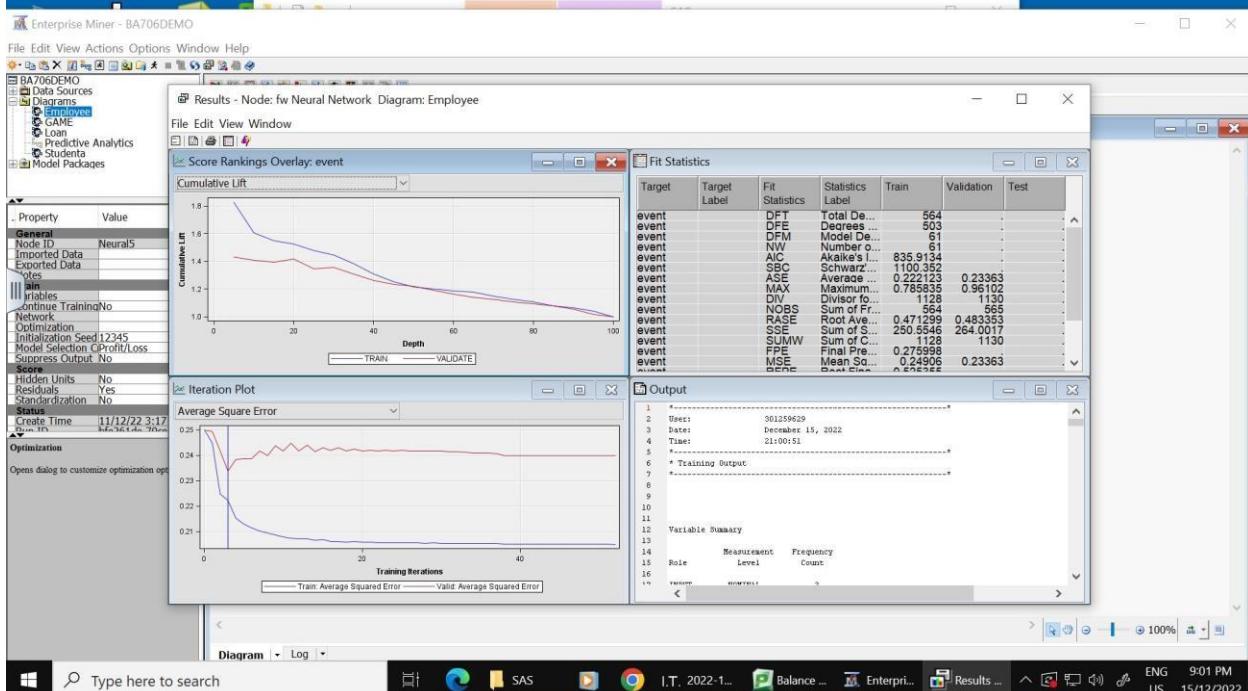
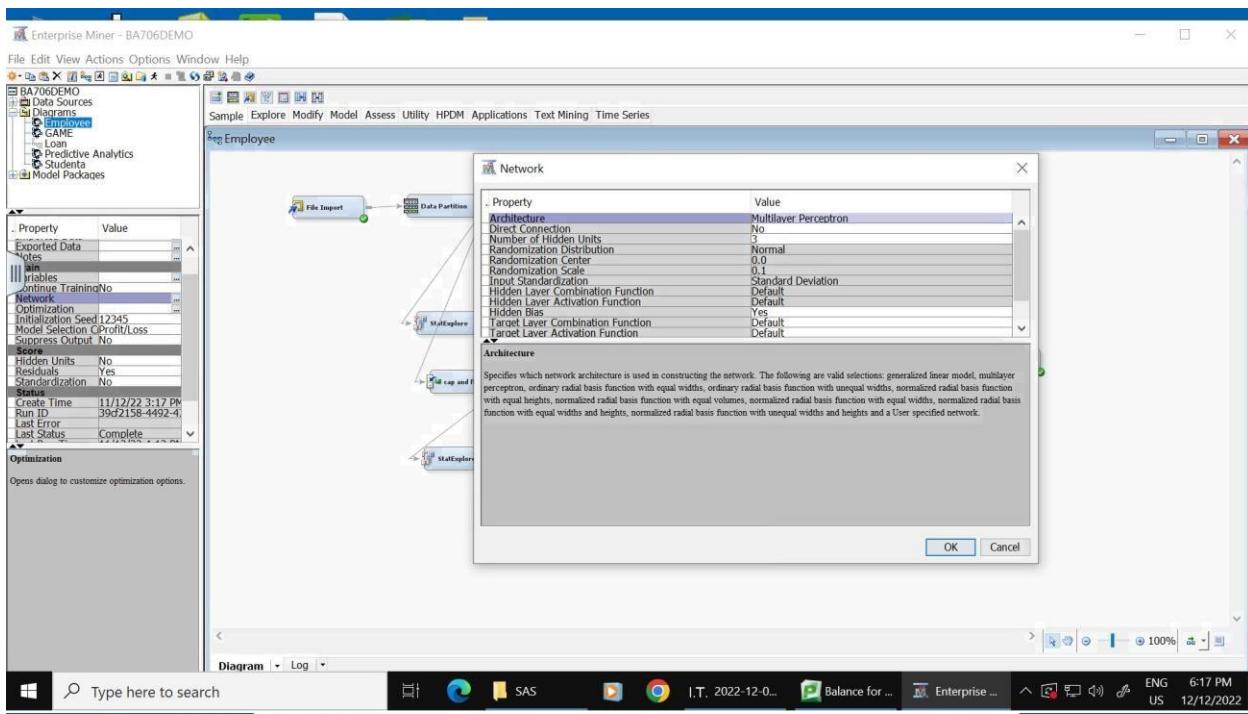
Here, we see that the average squared error of the training and validation data sets start off in opposite directions with the validation getting worse before it changes direction for improvement. The most improvement for the validation ASE was at iteration 6 then a slight increase occurred which is a sign of overfitting. At iteration 53 convergence occurred.



NEURAL NETWORK WITH 3 HIDDEN UNITS

This neural network node was run using three hidden units. The ASE was 0.23363 which is an improvement from the neural network with 2 hidden units.





The screenshot shows the Enterprise Miner interface with the following details:

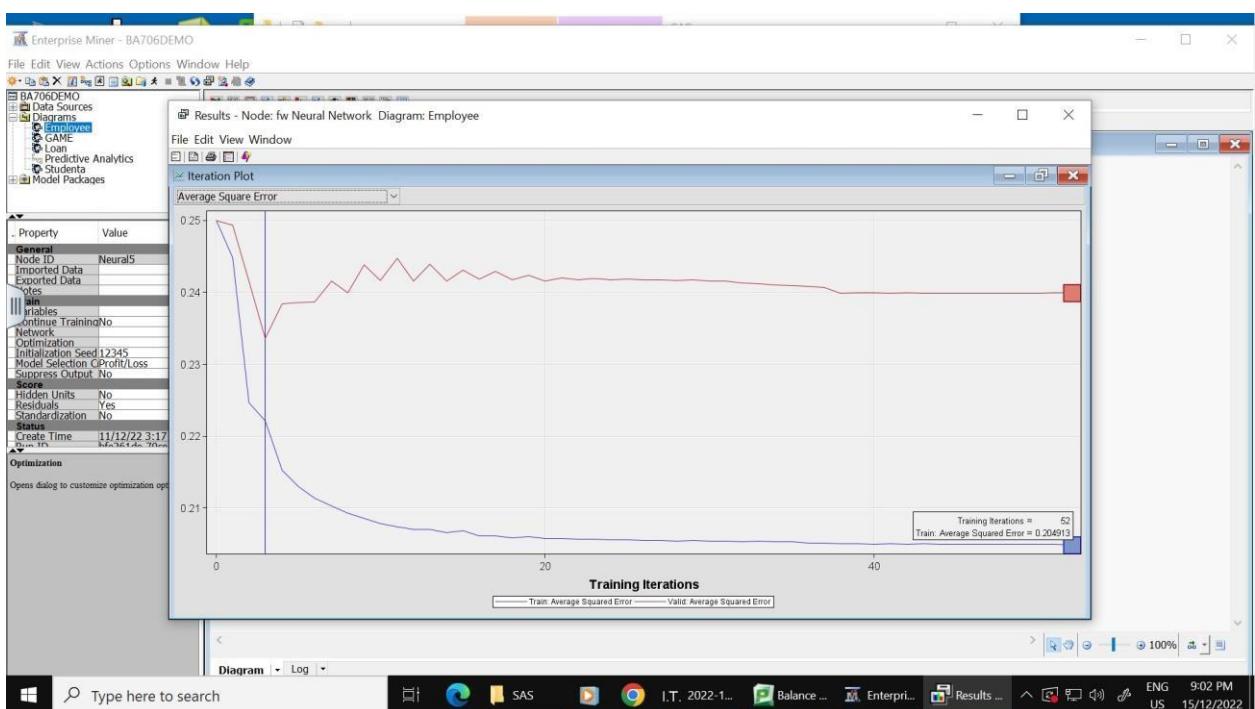
- Fit Statistics Window:**

Target	Target Label	Fit Statistics	Statistics Label	Train	Validation	Test
event		DFT	Total Degrees of Free...	564	...	
event		DFE	Degrees of Freedom f...	503	...	
event		DPM	Model Degree of Free...	61	...	
event		NW	Number of Estimated...	61	...	
event		AIC	Akaike's Information Cr...	835.9134	...	
event		SBC	Schwarz's Bayesian Cr...	1100.352	...	
event		ASE	Average Squared Error	0.23363	0.23363	0.23363
event		MAX	Maximum Absolute Error	0.785835	0.96102	0.96102
event		DIV	Divisor for ASE	1128	1130	1130
event		NBDS	Sum of Frequencies	364	358	358
event		RASE	Root Average Squared...	0.471299	0.483353	0.483353
event		SSE	Sum of Squared Errors	250.5546	264.0017	264.0017
event		SUMW	Sum of Case Weights ...	1128	1130	1130
event		FPE	Final Prediction Error	0.27488	0.24806	0.23363
event		MSE	Mean Squared Error	0.24806	0.23363	0.23363
event		RFPE	Root Final Prediction ...	0.523355	0.499059	0.483353
event		RMSE	Root Mean Squared E...	0.499059	0.483353	0.483353
event		AVERR	Average Error Function	0.683402	0.683402	0.683402
event		ERR	Error Function	713.9134	746.6914	746.6914
event		MISC	Misclassification Rate	0.388298	0.39646	0.39646
event		WRONG	Number of Wrong Clas...	219	224	224
- Iteration Plot Window:**

Average Square Error

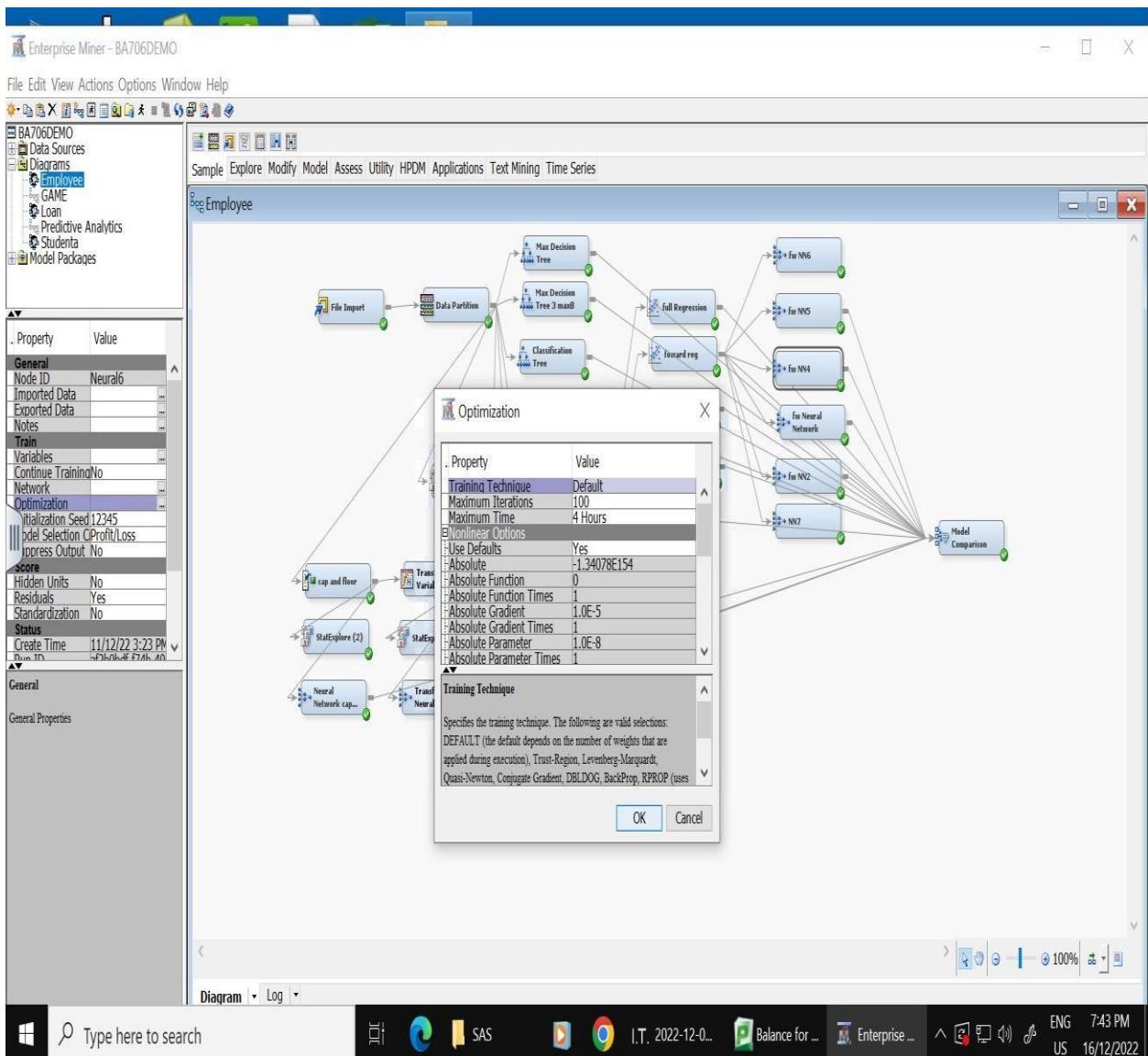
Training Iterations	Train: Average Squared Error	Valid: Average Squared Error
0	0.25	0.25
5	0.22	0.24
10	0.205	0.24
20	0.204	0.24
30	0.204	0.245
40	0.204	0.245
52	0.204	0.245

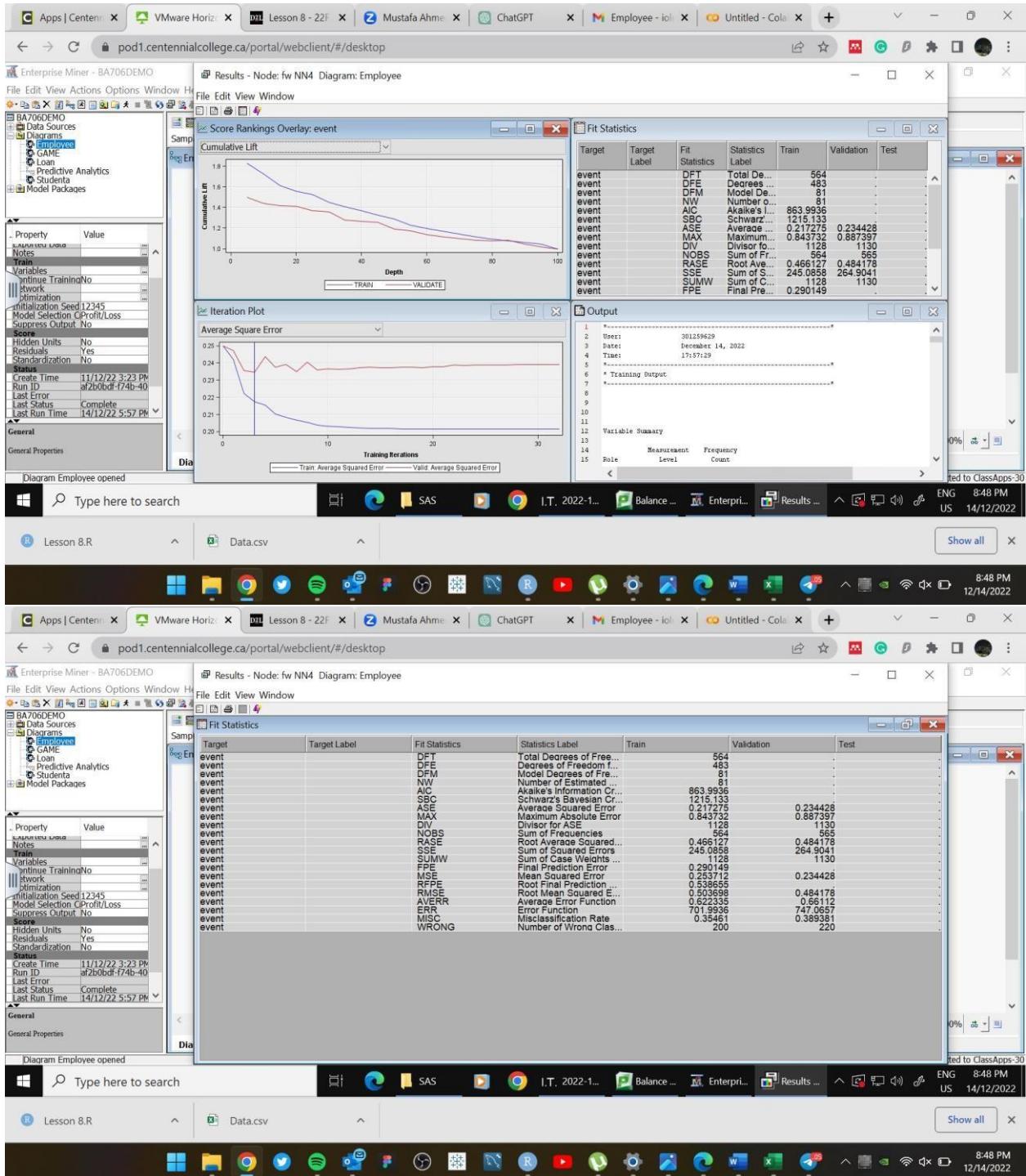
Here, we see that the average squared error of the training and validation data sets start off in the same direction of improvement till iteration 3 where the validation ASE starts to increase as a sign of overfitting. this occurs until the point of convergence at iteration 52.



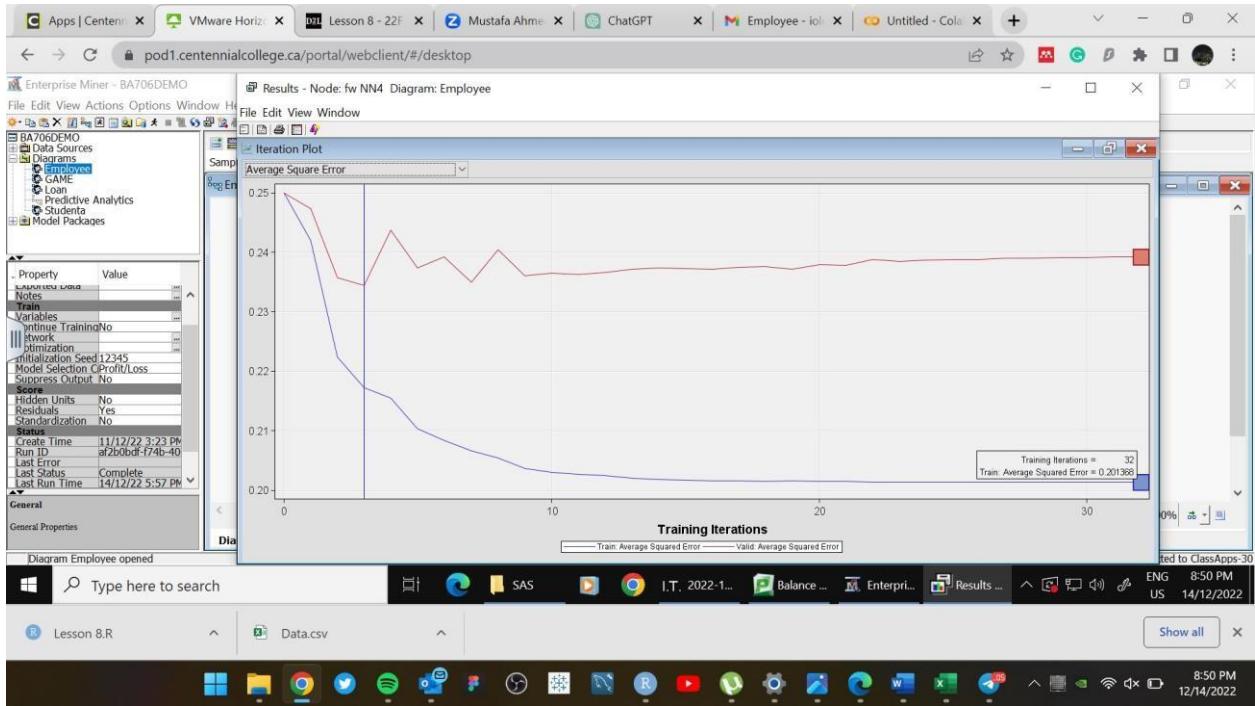
NEURAL NETWORK WITH 4 HIDDEN UNITS

We continued to check for the best neural network model by increasing the number of hidden units to 4 as we noticed 3 hidden units was an improvement when compared to 2 hidden units. The result from the four hidden units neural network was even worse than that of the two hidden units with an ASE of 0.234428 .



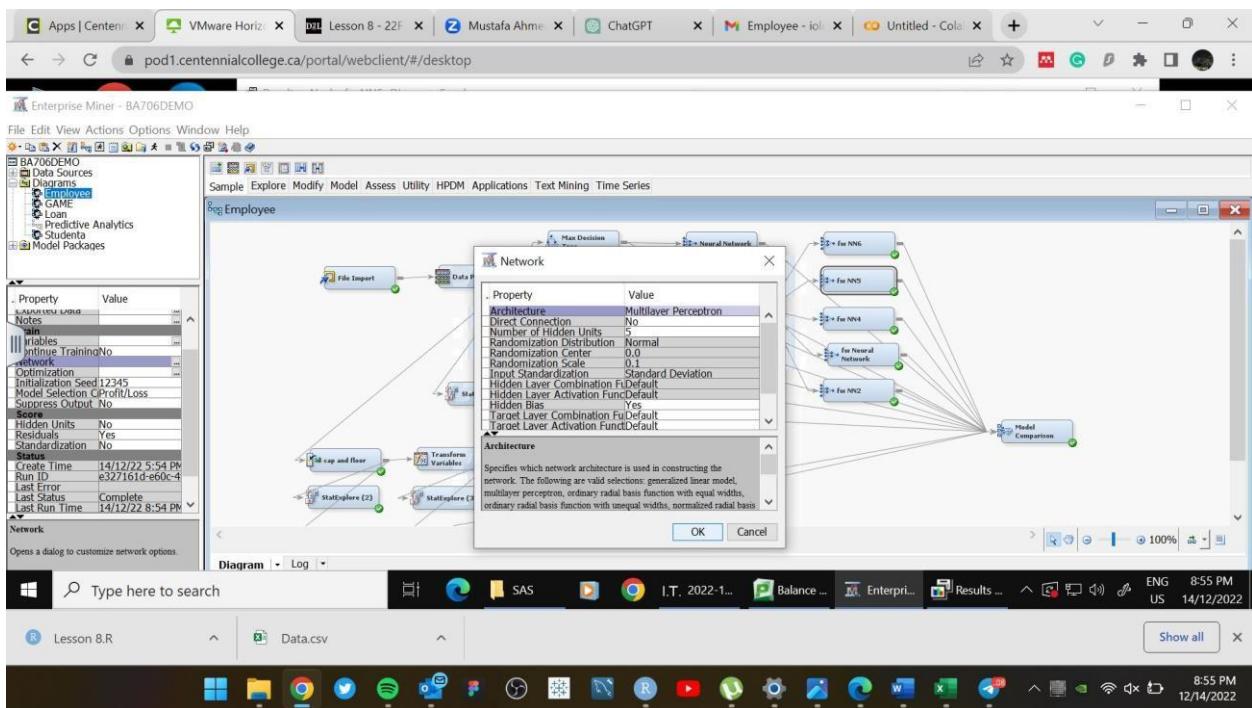
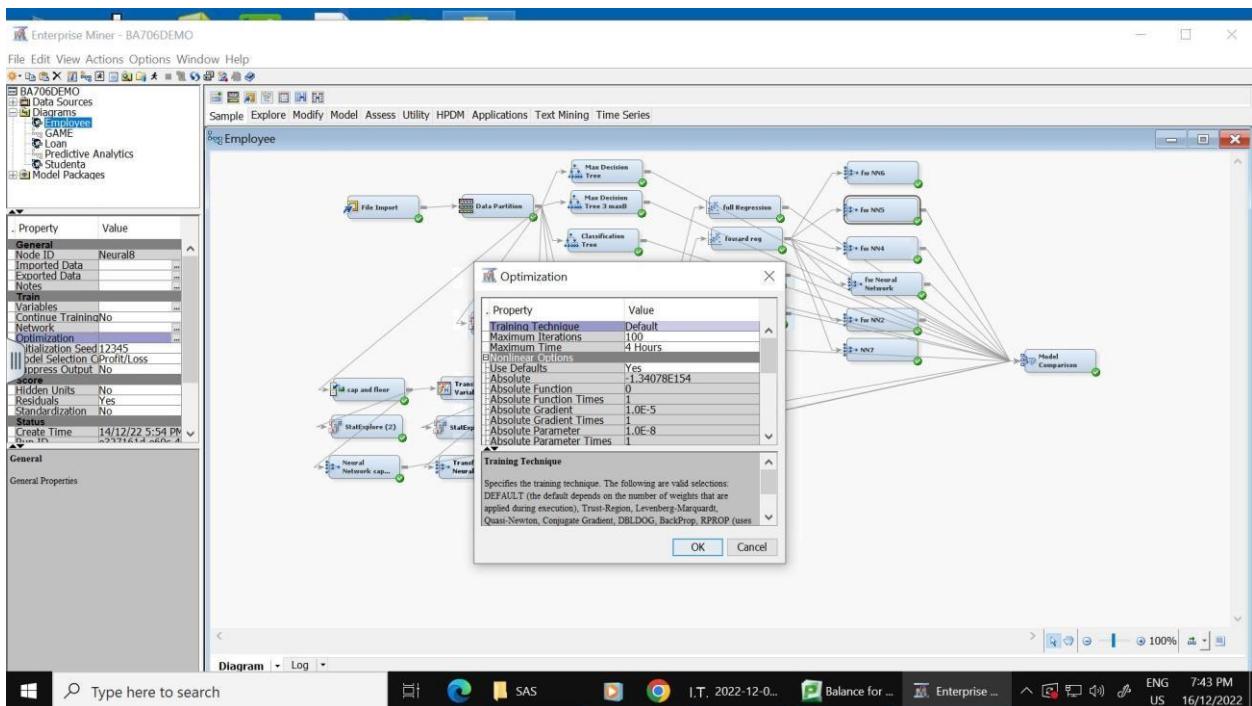


Here, we see that the average squared error of the training and validation data sets start off in the same direction of improvement and then the validation ASE starts to increase as a sign of overfitting. This occurs until the point of convergence at iteration 32.



NEURAL NETWORK WITH 5 HIDDEN UNITS

We further increased the number of hidden units to 5 and the outcome ASE of 0.234067 was an improvement on the two hidden unit neural network but not as good as the three hidden unit neural network. This proves no further need to run additional hidden units as the aim is to find the least complex model with the best result.



Enterprise Miner - BA706DEMO

File Edit View Actions Options Window Help

Score Rankings Overlay: event

Cumulative Lift

Fit Statistics

Target	Target Label	Fit Statistics	Statistics Label	Train	Validation	Test
event	DFT	Total Degrees of Free...	564			
event	DFE	Degrees of Freedom F...	463			
event	DFM	Model Degrees of Free...	101			
event	NW	Number of Estimated ...	101			
event	AIC	Akaike's Information Cr...	912.8829			
event	SSC	Schwarz's Inf...	1350.723			
event	ASE	Average Squared Error	0.220765	0.234067		
event	MAX	Maximum Absolute Error	0.844101	0.875518		
event	DIV	Divisor for ASE	1128	1130		
event	NBDS	Sum of Case Weights ...	564	565		
event	RASE	Root Average Squared E...	0.469856	0.483804		
event	SSE	Sum of Squared Errors	249.0224	264.4952		
event	SUMW	Sum of Case Weights ...	1128	1130		
event	FPE	Firth's Prediction Erro...	0.317081			
event	MSE	Mean Squared Error	0.268923	0.234067		
event	RFPE	Root Final Prediction E...	0.563099	0.483804		
event	RNSSE	Root Normalized Squared E...	0.513013	0.483804		
event	AVERR	Average Error Function	0.530215	0.660831		
event	ERR	Error Function	710.8829	746.7386		
event	MISC	Misclassification Rate	0.363475	0.40708		
event	WRONG	Number of Wrong Classi...	205	230		

Output

```

1 * User: 301159629
2 Date: 14/12/2022
3 Time: 20:54:53
4
5 * Training Output
6
7
8
9
10
11 Variable Summary
12
13
14 Role Measurement Level Frequency Count
< >

```

Lesson 8.R Data.csv

8:55 PM 12/14/2022

Enterprise Miner - BA706DEMO

File Edit View Actions Options Window Help

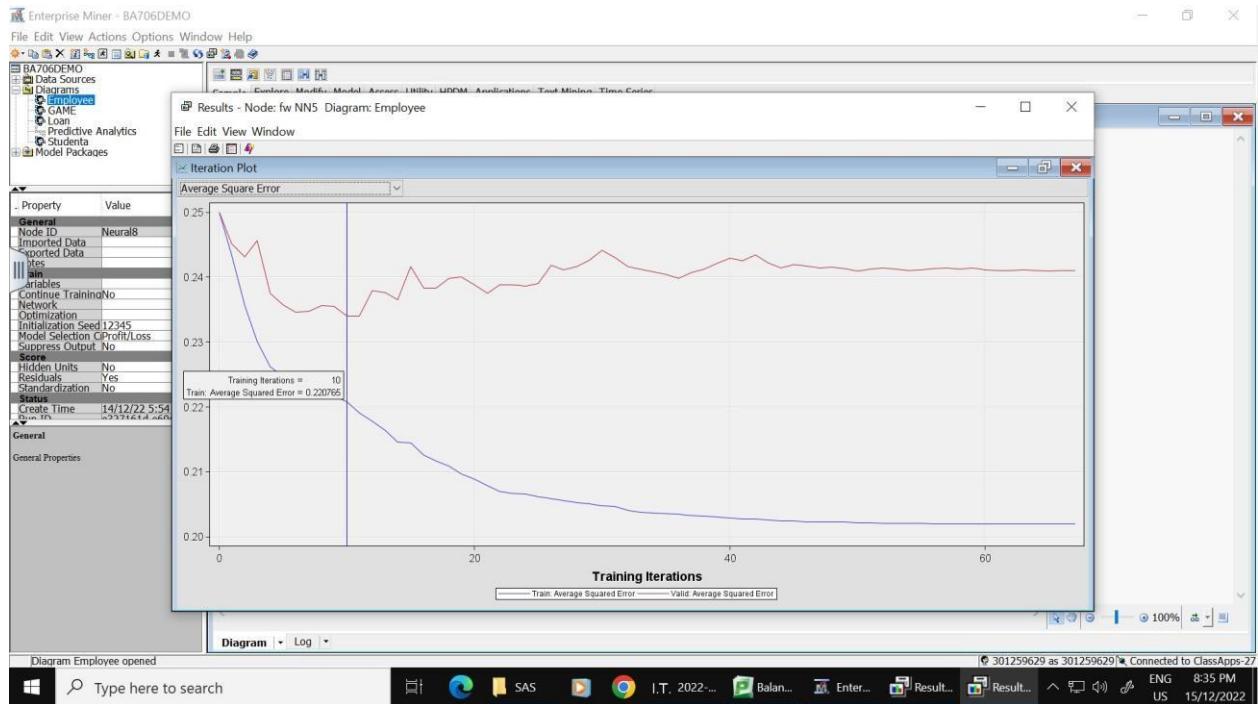
Fit Statistics

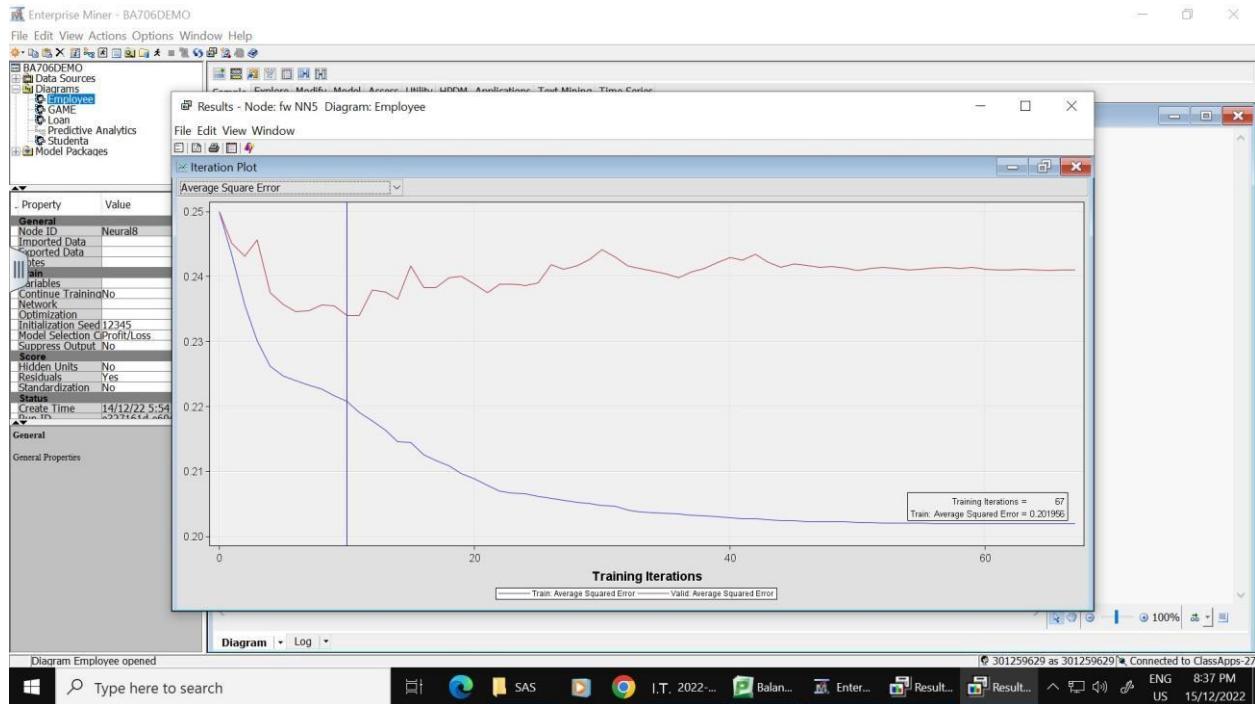
Target	Target Label	Fit Statistics	Statistics Label	Train	Validation	Test
event	DFT	Total Degrees of Free...	564			
event	DFE	Degrees of Freedom F...	463			
event	DFM	Model Degrees of Free...	101			
event	NW	Number of Estimated ...	101			
event	AIC	Akaike's Information Cr...	912.8829			
event	SSC	Schwarz's Inf...	1350.723			
event	ASE	Average Squared Error	0.220765	0.234067		
event	MAX	Maximum Absolute Error	0.844101	0.875518		
event	DIV	Divisor for ASE	1128	1130		
event	NBDS	Sum of Case Weights ...	564	565		
event	RASE	Root Average Squared E...	0.469856	0.483804		
event	SSE	Sum of Squared Errors	249.0224	264.4952		
event	SUMW	Sum of Case Weights ...	1128	1130		
event	FPE	Firth's Prediction Erro...	0.317081			
event	MSE	Mean Squared Error	0.268923	0.234067		
event	RFPE	Root Final Prediction E...	0.563099	0.483804		
event	RNSSE	Root Normalized Squared E...	0.513013	0.483804		
event	AVERR	Average Error Function	0.530215	0.660831		
event	ERR	Error Function	710.8829	746.7386		
event	MISC	Misclassification Rate	0.363475	0.40708		
event	WRONG	Number of Wrong Classi...	205	230		

Lesson 8.R Data.csv

8:56 PM 12/14/2022

Here, we see that the average squared error of the training and validation data sets start off in the same direction of improvement till iteration 10 where the validation ASE starts to increase as a sign of overfitting. This occurs until the point of convergence at iteration 67.

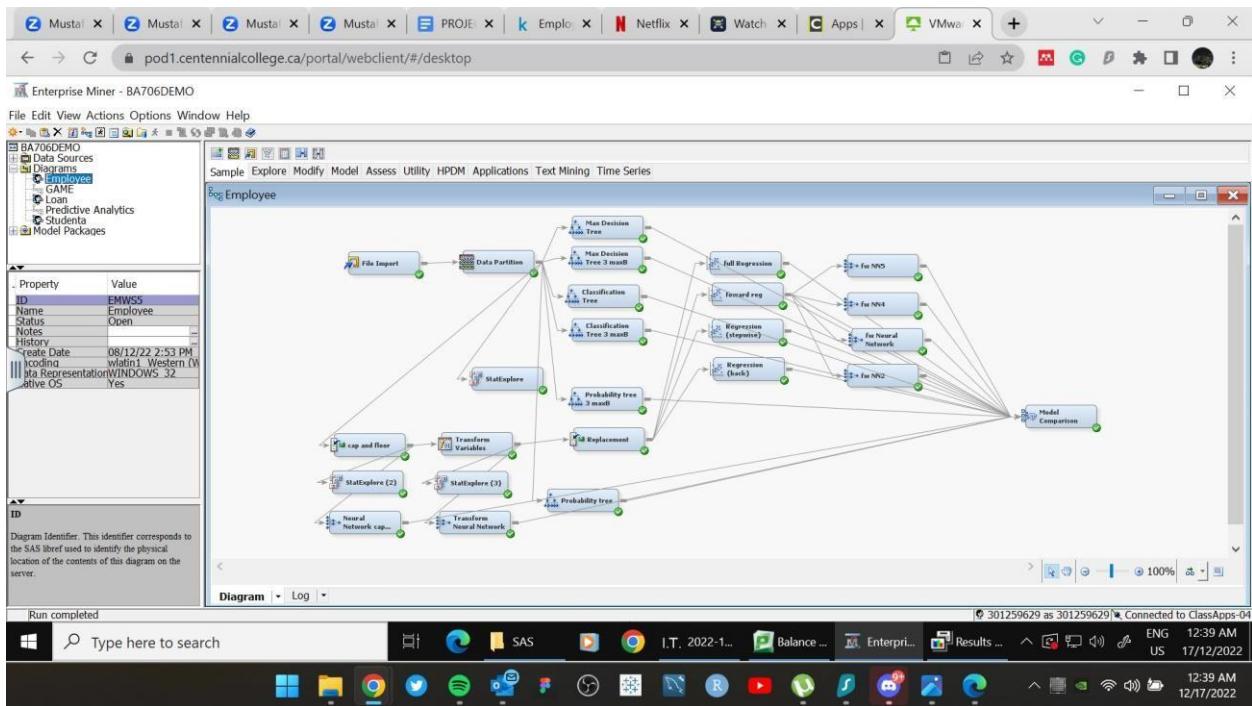




ASSESSMENT

MODEL COMPARISON

Finally, a model comparison node was introduced to the diagram workspace to which the different decision trees, regressions, and neural network models were connected. The results show a comparison of all models using several statistical tests. The below shows the final diagram of all the models we tested in this project:



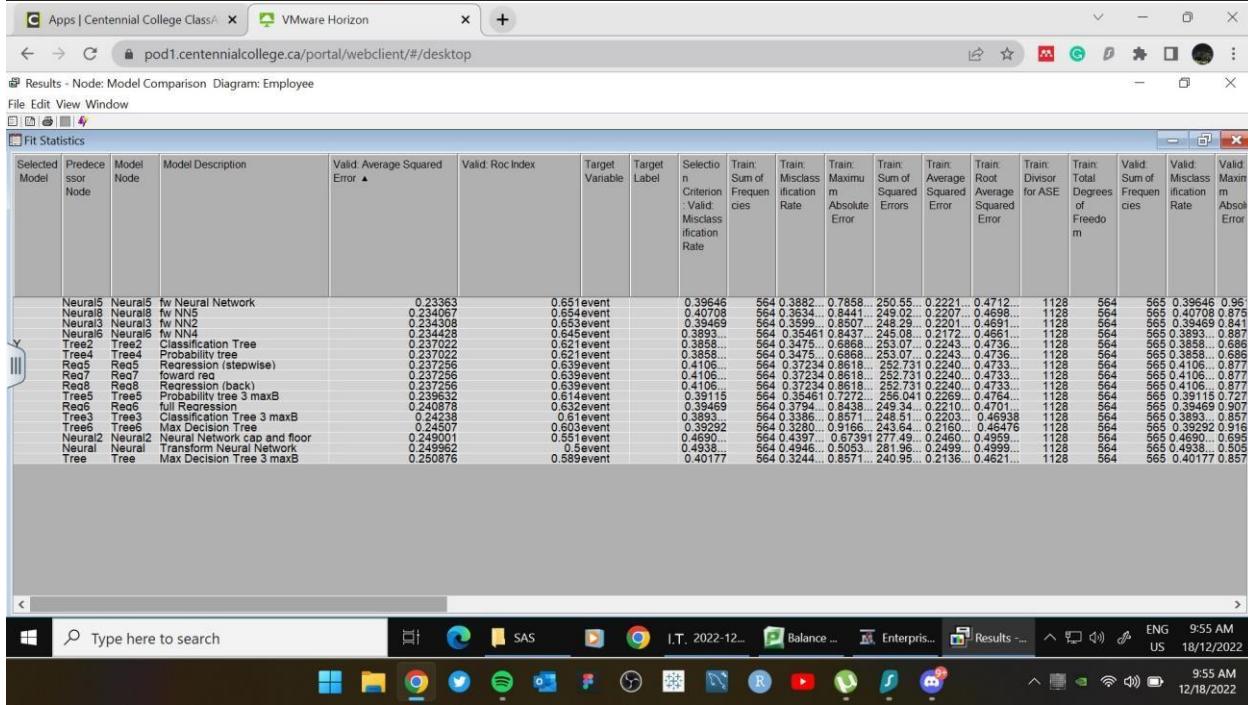
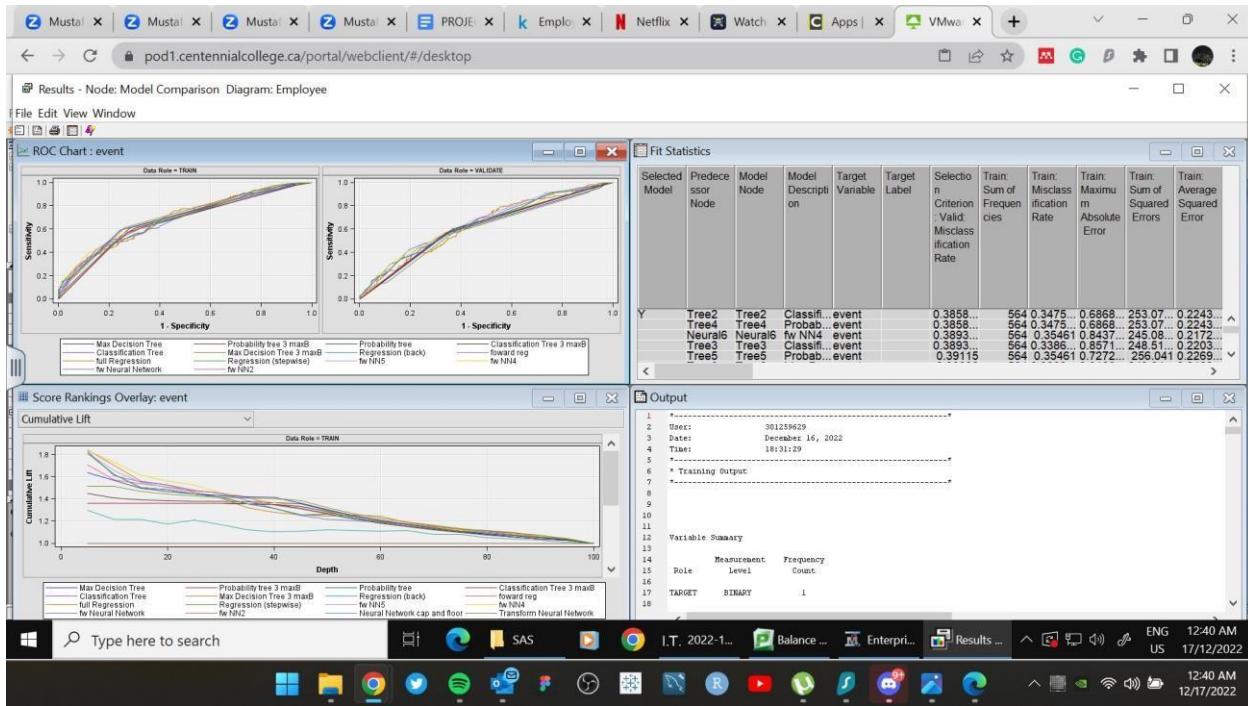
As stated earlier, we will be selecting the most fit model using the Average squared error. Below shows the summary of all models the Average squared error and ROC index. Using the ROC index as a selection criterion, the best model would be the Neural Network with 5 hidden units and ROC index of 0.654. However, our chosen method of selection is the Average squared error hence the Neural Network with 3 hidden units is the best model.

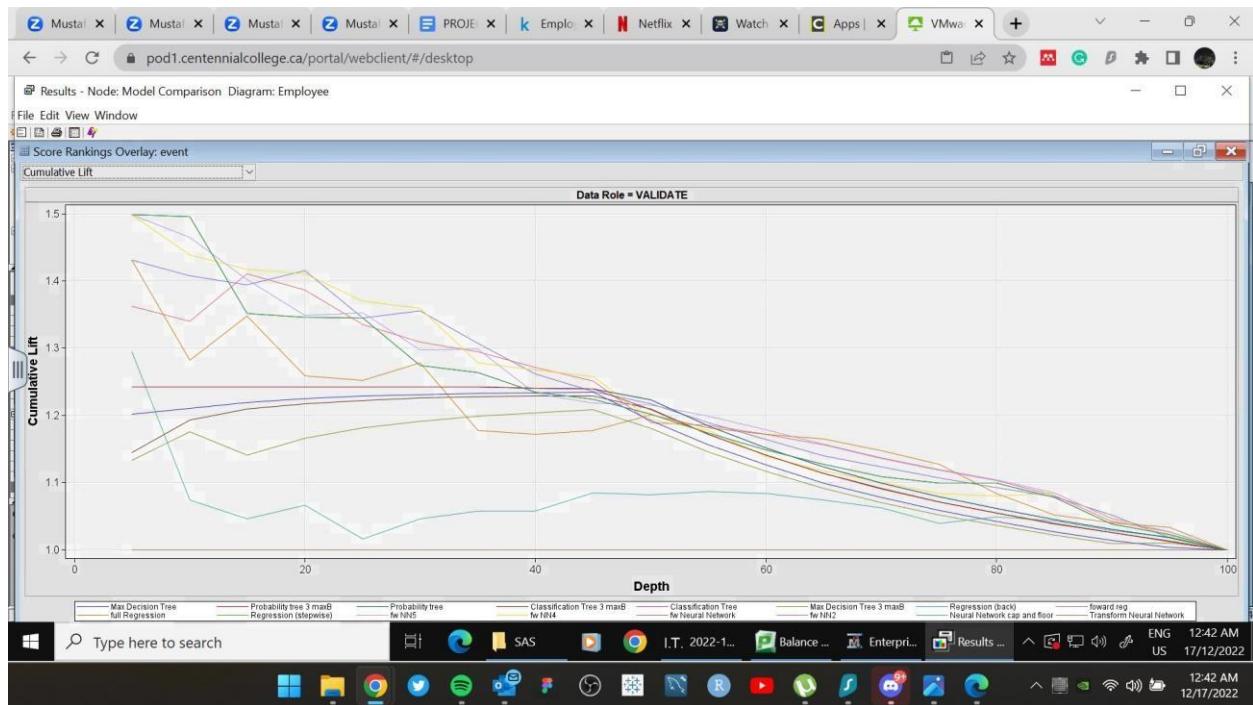
Model Type	Average Squared Error	ROC
------------	-----------------------	-----

Neural Network 3 hidden 0.233630 0.651 units

Neural Network 5 hidden units 0.234067 0.654

	0.234308	Neural	0.653
Network 2 hidden units			
	0.234428	Neural	0.645
Network 4 hidden units			
Classification Tree 2 branch	0.237022		0.621
Probability Tree 2 branch	0.237022		0.621
Regression Stepwise	0.237256		0.639
Regression Forward	0.237256		0.639
Regression Backward	0.237256		0.639
Probability Tree 3 branch	0.239632		0.614
Regression Full	0.240878		0.632
Classification Tree 3 branch	0.242380		0.610
Maximal Decision Tree	0.245070		0.603





CONCLUSION

The best model is the Neural Network with 3 hidden units which has a validation Average Square Error of 0.23363. This neural network model was derived from attaching a neural network node to the forward regression model which had the best ASE of 0.237256. However, because of the complexity of how the Neural network arrives at its result, we are unable to explain in detail the next steps to take to reduce attrition in the different industries using this model. For this reason, we will be using the next best model which is easy to interpret to make suggestions for improvement. From the comparison table, the next best models which are not neural network models are the classification tree and probability tree with two branches. They both have a Validation Average Square Error of 0.237022. Therefore, we will make our recommendations based on these trees. We will also use the forward regression model results to support our recommendations.

For all the trees used in this project, the first split occurred in the industry variable. This variable had the highest log worth and was considered the most important variable by the decision trees and interestingly by most of the regression models too. Other variables considered by the trees are Profession, Stag (duration of experience), and Way (means of transportation to the workplace).

The trees show that 60% of employees in Banks, Consult and State industries are likely to quit their jobs with 62.89% of these employees requiring a vehicle to get to work. Employees who walk to the office (most likely because they stay in close proximity) are less likely to quit their jobs. 43.99% of the employees in the power generation, and retail industries are likely to quit with a higher percentage than those in the commercial, sales, and finance professions.

The forward regression model used the transport industry as a means of comparison for all other industries. It appears to be the most stable industry that is able to retain its employees after the IT industry with the least likelihood of attrition. Employees in HR and IT profession when compared with sales are also less likely to quit their jobs.

Below are the insights obtained from the selected models and our recommendations:

1. **Banks, Consult, and State Industry:** With 60% of employees leaving their jobs, a deep dive could show specific reasons for attrition. From the forward regression model, those in

banks, consult, and state are 4.6, 4.4 and 3.1 times respectively more likely to quit their jobs when compared with those in the transport industry. There might also be a need to study the transport industry to see why it seems more stable and able to retain its employees.

We suggest a drive for work-life balance is implemented in these industries. Incentives should be given for working longer than usual work hours, compensation given for achieving goals far and beyond expectations, and recognition for assignments well executed. A program for continued benefits after retirement should also be introduced. All these might be useful in encouraging professionals within these industries to stay on the job for a longer period or till retirement.

2. **Employees who come to work by bus or car:** Proximity to the office location also appears to be an important factor from both the tree models and regression models. People who have to get to work with a vehicle are 2.5 times more likely to quit their jobs. We advise that attention is given to how employees get to work. A transportation initiative can be introduced whereby employees have a choice of joining a staff bus that does a pick-up and drop-off for staff who stay a distance from the workplace. This will reduce the stress of driving to work or joining the public transportation system on a daily basis. In addition, the employees can be given the option of working remotely on some days in the week. This option has become a widely accepted work arrangement in the past two years which has improved the work-life balance of employees.

3. **Commercial, Sales and Finance Profession:** From the regression models, we see that employees in all other professions including sales are more likely to quit their jobs when compared with HR and IT professionals. We would say carrying out an investigation on the two professions might give an insight as to the reason for their stability. In the mean time, suggestions made for the industries above will also be profitable when applied to the profession variable. In addition, we suggest significance should be attached to attaining specific levels of achievement. This can be an incentive for employees to spend more time on their jobs while they pursue higher levels in their career paths.

REFERENCE

<https://www.kaggle.com/datasets/davinwijaya/employee-turnover>